

# Nonlinear differential equation problems

Hans Petter Langtangen

Medical Computing, Simula Research Laboratory and Department of Informatics, University of Oslo

Nov 5, 2014

Note: This note is still under development.

## Contents

<b>Introduction of basic concepts</b>	<b>3</b>
1.1 Linearization by explicit time discretization . . . . .	3
1.2 Exact solution of nonlinear equations . . . . .	4
1.3 Linearization . . . . .	5
1.4 Picard iteration . . . . .	5
1.5 Linearization by a geometric mean . . . . .	7
1.6 Newton's method . . . . .	8
1.7 Relaxation . . . . .	10
1.8 Implementation and experiments . . . . .	10
1.9 Generalization to a general nonlinear ODE . . . . .	14
<b>Systems of nonlinear algebraic equations</b>	<b>15</b>
2.1 Picard iteration . . . . .	15
2.2 Newton's method . . . . .	16
2.3 Stopping criteria . . . . .	18
2.4 Example: A nonlinear ODE model from epidemiology . . . . .	19
<b>Linearization at the differential equation level</b>	<b>20</b>
3.1 Explicit time integration . . . . .	21
3.2 Backward Euler scheme and Picard iteration . . . . .	21
3.3 Backward Euler scheme and Newton's method . . . . .	22
3.4 Crank-Nicolson discretization . . . . .	24
<b>Discretization of stationary nonlinear differential equations</b>	<b>24</b>
4.1 Finite difference discretizations . . . . .	25
4.2 Solution of algebraic equations . . . . .	26
4.3 Galerkin-type discretizations . . . . .	28
4.4 Finite element basis functions . . . . .	29
4.5 The group finite element method . . . . .	30

4.6 Numerical integration of nonlinear terms . . . . .	
4.7 Finite element discretization of a variable coefficient Laplace equation . . . . .	
4.8 Picard iteration defined from the variational form . . . . .	
4.9 Newton's method defined from the variational form . . . . .	

## 5 Multi-dimensional PDE problems

5.1 Finite element discretization . . . . .	
5.2 Finite difference discretization . . . . .	
5.3 Continuation methods . . . . .	

## 6 Exercises

In a linear differential equation all terms involving the unknown function are linear in the unknown functions or their derivatives. Linear here means that the unknown function or a derivative of it is multiplied by a number or a known function. All other differential equations are non-linear. The easiest way to see if an equation is nonlinear is to spot nonlinear terms where the unknown function or their derivatives are multiplied by each other. For example, in

$$u'(t) = -a(t)u(t) + b(t),$$

the terms involving the unknown function  $u$  are linear:  $u'$  contains the derivative of the unknown function multiplied by unity, and  $au$  contains the unknown function multiplied by a known function. However,

$$u'(t) = u(t)(1 - u(t)),$$

the equation is nonlinear because of the term  $-u^2$  where the unknown function is multiplied by itself. Also

$$\frac{\partial u}{\partial t} + u \frac{\partial u}{\partial x} = 0,$$

the equation is nonlinear because of the term  $uu_x$  where the unknown function appears as a product with itself or one of its derivatives. Another example of a nonlinear equation is

$$u'' + \sin(u) = 0,$$

because  $\sin(u)$  contains products of  $u$ ,

$$\sin(u) = u - \frac{1}{3}u^3 + \dots$$

A series of forthcoming examples will explain how to tackle nonlinear differential equations with various techniques.

## Introduction of basic concepts

Consider the (scaled) logistic equation

$$u'(t) = u(t)(1 - u(t)). \quad (1)$$

This is a nonlinear differential equation which will be solved by different strategies in the following. A time discretization of (1) will either lead to a linear algebraic equation or a nonlinear algebraic equation at each time level. In the former case, the time discretization method transforms the nonlinear ODE into linear subproblems at each time level, and the solution is straightforward to find since linear algebraic equations are easy to solve by hand. However, when the time discretization leads to nonlinear algebraic equations, we cannot (except in very rare cases) solve these without turning to approximate, iterative solution methods.

The following subsections first introduce various methods using (1):

- explicit time discretization methods (with no need to solve nonlinear algebraic equations)
- implicit Backward Euler discretization, leading to nonlinear algebraic equations solved by
  - an exact analytical technique
  - Picard iteration based on manual linearization
  - a single Picard step
  - Newton's method
- Implicit Crank-Nicolson discretization and linearization via a geometric mean formula

Hereafter, we compare the performance of the various approaches. Despite the simplicity of (1), the conclusions reveal typical features of the various methods in much more complicated nonlinear PDE problems.

### 1.1 Linearization by explicit time discretization

Forward Euler method to solve (1) results in

$$\frac{u^{n+1} - u^n}{\Delta t} = u^n(1 - u^n),$$

which is a *linear* algebraic equation for the unknown value  $u^{n+1}$ . The nonlinearity in the original equation poses in this case no difficulty in the discrete algebraic equation. Any other explicit scheme in time will also give only linear algebraic equations to solve. For example, a typical 2nd-order Runge-Kutta method for (1) reads,

$$\begin{aligned} u^* &= u^n + \Delta t u^n(1 - u^n), \\ u^{n+1} &= u^n + \Delta t \frac{1}{2} (u^n(1 - u^n) + u^*(1 - u^*)) . \end{aligned}$$

The first step is linear in the unknown  $u^*$ . Then  $u^*$  is known in the next step, which is linear in the unknown  $u^{n+1}$ .

### 1.2 Exact solution of nonlinear equations

Switching to a Backward Euler scheme for (1),

$$\frac{u^n - u^{n-1}}{\Delta t} = u^n(1 - u^n),$$

this results in a nonlinear algebraic equation for the unknown value  $u^n$ . The equation is of quadratic type:

$$\Delta t(u^n)^2 + (1 - \Delta t)u^n - u^{n-1} = 0.$$

We shall now introduce a shorter and often cleaner notation for nonlinear algebraic equations at a given time level. The notation is inspired by the notation, i.e., variable names, used in a program, especially in more advanced partial differential equation problems. The unknown in the algebraic equation is denoted by  $u$ , while  $u^{(1)}$  is the value of the unknown at the previous time level (in general  $u^{(\ell)}$  is the value of the unknown  $\ell$  levels back in time). The notation will be frequently used in later sections. What is meant by  $u$  (the exact solution of the PDE problem, the numerical approximation to the exact solution or unknown solution at a certain time level) should be evident from the context.

The quadratic equation for the unknown  $u^n$  in (2) can with the new notation be written

$$F(u) = \Delta t u^2 + (1 - \Delta t)u - u^{(1)} = 0.$$

The solution is readily found to be

$$u = \frac{1}{2\Delta t} \left( -1 + \Delta t \pm \sqrt{(1 - \Delta t)^2 - 4\Delta t u^{(1)}} \right).$$

Now we encounter a fundamental challenge with nonlinear algebraic equations: the equation may have more than one solution. How do we pick the right solution? In the present simple case we can expand the square root in a series and truncate after the linear term since the Backward Euler scheme will in any case have an error proportional to  $\Delta t$  anyway. Using `sympy` we find the following series expansions of the roots:

```
>>> import sympy as sp
>>> dt, u_1, u = sp.symbols('dt u_1 u')
>>> r1, r2 = sp.solve(dt*u**2 + (1-dt)*u - u_1, u) # find roots
>>> r1
(dt - sqrt(dt**2 + 4*dt*u_1 - 2*dt + 1) - 1)/(2*dt)
>>> r2
(dt + sqrt(dt**2 + 4*dt*u_1 - 2*dt + 1) - 1)/(2*dt)
>>> print r1.series(dt, 0, 2)
-1/dt + 1 - u_1 + dt*(u_1**2 - u_1) + O(dt**2)
>>> print r2.series(dt, 0, 2)
1_1 + dt*(-u_1**2 + u_1) + O(dt**2)
```

We see that the `r1` root, corresponding to a minus sign in front of the square root in (4), behaves as  $1/\Delta t$  and will therefore blow up as  $\Delta t \rightarrow 0$ ! Therefore, only the `r2` root is of relevance in this case.

### .3 Linearization

When the time integration of an ODE results in a nonlinear algebraic equation, we must normally find its solution by defining a sequence of linear equations and hope that the solutions of these linear equations converge to the desired solution of the nonlinear algebraic equation. Usually this means solving the linear equation repeatedly in an iterative fashion. Alternatively, the nonlinear equation can sometimes be approximated by one linear equation, and consequently there is no need for iteration.

Constructing a linear equation from a nonlinear one requires *linearization* of each nonlinear term. This can be done manually as in Picard iteration, or implicitly algorithmically as in Newton's method. Examples will best illustrate how to linearize nonlinear problems.

### .4 Picard iteration

Let us write (3) in a more compact form

$$F(u) = au^2 + bu + c = 0,$$

with  $a = \Delta t$ ,  $b = 1 - \Delta t$ , and  $c = -u^{(1)}$ . Let  $u^-$  be an available approximation of the unknown  $u$ . Then we can linearize the term  $u^2$  simply by writing  $u^-u$ . The resulting equation,  $\hat{F}(u) = 0$ , is now linear and hence easy to solve:

$$F(u) \approx \hat{F}(u) = au^-u + bu + c = 0.$$

Since the equation  $\hat{F} = 0$  is only approximate, the solution  $u$  does not equal the exact solution  $u_e$  of the exact equation  $F(u_e) = 0$ , but we can hope that  $u$  is closer to  $u_e$  than  $u^-$  is, and hence it makes sense to repeat the procedure, i.e., at  $u^- = u$  and solve  $\hat{F}(u) = 0$  again.

The idea of turning a nonlinear equation into a linear one by using an approximation  $u^-$  of  $u$  in nonlinear terms is a widely used approach that goes under many names: *fixed-point iteration*, the method of *successive substitutions*,

*nonlinear Richardson iteration*, and *Picard iteration*. We will stick to the latter name.

Picard iteration for solving the nonlinear equation arising from the backward Euler discretization of the logistic equation can be written as

$$u = -\frac{c}{au^- + b}, \quad u^- \leftarrow u.$$

The iteration is started with the value of the unknown at the previous time level  $u^- = u^{(1)}$ .

Some prefer an explicit iteration counter as superscript in the mathematical notation. Let  $u^k$  be the computed approximation to the solution in iteration  $k$ . In iteration  $k + 1$  we want to solve

$$au^k u^{k+1} + bu^{k+1} + c = 0 \quad \Rightarrow \quad u^{k+1} = -\frac{c}{au^k + b}, \quad k = 0, 1, \dots$$

Since we need to perform the iteration at every time level, the time level index is often also included:

$$au^{n+1,k} u^{n+1,k+1} + bu^{n+1,k+1} - u^n = 0 \quad \Rightarrow \quad u^{n+1,k+1} = \frac{u^n}{au^{n+1,k} + b},$$

with the start value  $u^{n,0} = u^{n-1}$  and the final converged value  $u^{n+1} = u^{n+1,k}$  for sufficiently large  $k$ .

However, we will normally apply a mathematical notation in our final result that is as close as possible to what we aim to write in a computer code. Then it becomes natural to use  $u$  and  $u^-$  instead of  $u^{k+1}$  and  $u^k$  or  $u^{n+1,k}$  and  $u^{n+1,k}$ .

**Stopping criteria.** The iteration method can typically be terminated when the change in the solution is smaller than a tolerance  $\epsilon_u$ :

$$|u - u^-| \leq \epsilon_u,$$

or when the residual in the equation is sufficiently small ( $\epsilon_r$ ),

$$|F(u)| = |au^2 + bu + c| < \epsilon_r.$$

**A single Picard iteration.** Instead of iterating until a stopping criterion is fulfilled, one may iterate a specific number of times. Just one Picard iteration is popular as this corresponds to the intuitive idea of approximating a nonlinear term like  $(u^n)^2$  by  $u^{n-1}u^n$ . This follows from the linearization  $u^-u^n$  with an initial choice of  $u^- = u^{n-1}$  at time level  $t_n$ . In other words, a single iteration corresponds to using the solution at the previous time level to linearize the nonlinear terms. The resulting discretization becomes

$$\frac{u^n - u^{n-1}}{\Delta t} = u^n(1 - u^{n-1}), \quad (5)$$

which is a linear algebraic equation in the unknown  $u^n$ , and therefore we can easily solve for  $u^n$ , and there is no need for any alternative notation.

We shall later refer to the strategy of taking one Picard step, or equivalently, linearizing terms with use of the solution at the previous time step, as the *Picard1* method. It is a widely used approach in science and technology, but with some limitations if  $\Delta t$  is not sufficiently small (as will be illustrated later).

#### Notice.

Equation (5) does not correspond to a “pure” finite difference method where the equation is sampled at a point and derivatives replaced by differences (because the  $u^{n-1}$  term on the right-hand side must then be  $u^n$ ). The best interpretation of the scheme (5) is a Backward Euler difference combined with a single (perhaps insufficient) Picard iteration at each time level, with the value at the previous time level as start for the Picard iteration.

## 5 Linearization by a geometric mean

We consider now a Crank-Nicolson discretization of (1). This means that the time derivative is approximated by a centered difference,

$$[D_t u = u(1 - u)]^{n+\frac{1}{2}},$$

written out as

$$\frac{u^{n+1} - u^n}{\Delta t} = u^{n+\frac{1}{2}} - (u^{n+\frac{1}{2}})^2. \quad (6)$$

The term  $u^{n+\frac{1}{2}}$  is normally approximated by an arithmetic mean,

$$u^{n+\frac{1}{2}} \approx \frac{1}{2}(u^n + u^{n+1}),$$

which means that the scheme involves the unknown function only at the time levels where we actually compute it. The same arithmetic mean applied to the nonlinear term gives

$$(u^{n+\frac{1}{2}})^2 \approx \frac{1}{4}(u^n + u^{n+1})^2,$$

which is nonlinear in the unknown  $u^{n+1}$ . However, using a *geometric mean* for  $(u^{n+\frac{1}{2}})^2$  is a way of linearizing the nonlinear term in (6):

$$(u^{n+\frac{1}{2}})^2 \approx u^n u^{n+1}.$$

Using an arithmetic mean on the linear  $u^{n+\frac{1}{2}}$  term in (6) and a geometric mean for the second term, results in a linearized equation for the unknown  $u^{n+1}$

$$\frac{u^{n+1} - u^n}{\Delta t} = \frac{1}{2}(u^n + u^{n+1}) + u^n u^{n+1},$$

which can readily be solved:

$$u^{n+1} = \frac{1 + \frac{1}{2}\Delta t}{1 + \Delta t u^n - \frac{1}{2}\Delta t} u^n.$$

This scheme can be coded directly, and since there is no nonlinear algebraic equation to iterate over, we skip the simplified notation with  $u$  for  $u^n$  and  $u^{(1)}$  for  $u^{n+1}$ . The technique with using a geometric average is an exact transformation of a nonlinear algebraic equation to a linear one, without a need for iterations.

The geometric mean approximation is often very effective for linearizing quadratic nonlinearities. Both the arithmetic and geometric mean approximations have truncation errors of order  $\Delta t^2$  and are therefore compatible with the truncation error  $\mathcal{O}(\Delta t^3)$  of the centered difference approximation for the Crank-Nicolson method.

Applying the operator notation for the means and finite differences, the linearized Crank-Nicolson scheme for the logistic equation can be compactly expressed as

$$[D_t u = \bar{u}^t + \overline{u^2}^{t,g}]^{n+\frac{1}{2}}.$$

#### Remark.

If we use an arithmetic instead of a geometric mean for the nonlinear term in (6), we end up with a nonlinear term  $(u^{n+\frac{1}{2}})^2$ . This term can be linearized as  $u^n u^{n+1}$  in a Picard iteration approach and in particular as  $u^n u^{n+1}$  in a Picard1 iteration approach. The latter gives a scheme almost identical to the one arising from a geometric mean (the difference in being  $\frac{1}{4}\Delta t u^n(u^{n+1} - u^n) \approx \frac{1}{4}\Delta t^2 u' u$ , i.e., a difference of  $\mathcal{O}(\Delta t^2)$ ).

## 1.6 Newton’s method

The Backward Euler scheme (2) for the logistic equation leads to a nonlinear algebraic equation (3). Now we write any nonlinear algebraic equation in its general and compact form

$$F(u) = 0.$$

Newton’s method linearizes this equation by approximating  $F(u)$  by its series expansion around a computed value  $u^-$  and keeping only the linear

$$F(u) = F(u^-) + F'(u^-)(u - u^-) + \frac{1}{2}F''(u^-)(u - u^-)^2 + \dots \\ \approx F(u^-) + F'(u^-)(u - u^-) = \hat{F}(u).$$

The linear equation  $\hat{F}(u) = 0$  has the solution

$$u = u^- - \frac{F(u^-)}{F'(u^-)}.$$

Expressed with an iteration index in the unknown, Newton's method takes on the more familiar mathematical form

$$u^{k+1} = u^k - \frac{F(u^k)}{F'(u^k)}, \quad k = 0, 1, \dots$$

It can be shown that the error in iteration  $k + 1$  of Newton's method is the square of the error in iteration  $k$ , a result referred to as *quadratic convergence*. This means that for small errors the method converges very fast, and in particular much faster than Picard iteration and other iteration methods. (The proof of this result is found in most textbooks on numerical analysis.) However, the quadratic convergence appears only if  $u^k$  is sufficiently close to the solution. Further away from the solution the method can easily converge very slowly or diverge. The reader is encouraged to do Exercise 3 to get a better understanding of the behavior of the method.

Application of Newton's method to the logistic equation discretized by the backward Euler method is straightforward as we have

$$F(u) = au^2 + bu + c, \quad a = \Delta t, \quad b = 1 - \Delta t, \quad c = -u^{(1)},$$

and then

$$F'(u) = 2au + b.$$

The iteration method becomes

$$u = u^- + \frac{a(u^-)^2 + bu^- + c}{2au^- + b}, \quad u^- \leftarrow u. \quad (7)$$

At each time level, we start the iteration by setting  $u^- = u^{(1)}$ . Stopping criteria as listed for the Picard iteration can be used also for Newton's method.

An alternative mathematical form, where we write out  $a$ ,  $b$ , and  $c$ , and use a time level counter  $n$  and an iteration counter  $k$ , takes the form

$$u^{n,k+1} = u^{n,k} + \frac{\Delta t(u^{n,k})^2 + (1 - \Delta t)u^{n,k} - u^{n-1}}{2\Delta t u^{n,k} + 1 - \Delta t}, \quad u^{n,0} = u^{n-1}, \quad k = 0, 1, \dots \quad (8)$$

The program implementation is much closer to (7) than to (8), but the latter is better aligned with the established mathematical notation used in the literature.

## 1.7 Relaxation

One iteration in Newton's method or Picard iteration consists of solving the problem  $\hat{F}(u) = 0$ . Sometimes convergence problems arise because the solution  $u$  of  $\hat{F}(u) = 0$  is "too far away" from the previously computed  $u^-$ . A remedy is to introduce a relaxation, meaning that we first solve  $\hat{F}(u)$  for a suggested value  $u^*$  and then we take  $u$  as a weighted mean of what  $u^-$ , and what our linearized equation  $\hat{F} = 0$  suggests,  $u^*$ :

$$u = \omega u^* + (1 - \omega)u^-.$$

The parameter  $\omega$  is known as a *relaxation parameter*, and a choice  $\omega < 1$  prevents divergent iterations.

Relaxation in Newton's method can be directly incorporated in the iteration formula:

$$u = u^- - \omega \frac{F(u^-)}{F'(u^-)}.$$

## 1.8 Implementation and experiments

The program `logistic.py`<sup>1</sup> contains implementations of all the methods described above. Below is an extract of the file showing how the Picard and Newton methods are implemented for a Backward Euler discretization of the equation.

```
def BE_logistic(u0, dt, Nt, choice='Picard',
               eps_r=1E-3, omega=1, max_iter=1000):
    if choice == 'Picard1':
        choice = 'Picard'
        max_iter = 1

    u = np.zeros(Nt+1)
    iterations = []
    u[0] = u0
    for n in range(1, Nt+1):
        a = dt
        b = 1 - dt
        c = -u[n-1]

        if choice == 'Picard':

            def F(u):
                return a*u**2 + b*u + c

            u_ = u[n-1]
            k = 0
            while abs(F(u_)) > eps_r and k < max_iter:
                u_ = omega*(-c/(a*u_ + b)) + (1-omega)*u_
                k += 1
            u[n] = u_
            iterations.append(k)
```

<sup>1</sup><http://tinyurl.com/nm5587k/nonlin/logistic.py>

```

elif choice == 'Newton':

    def F(u):
        return a*u**2 + b*u + c

    def dF(u):
        return 2*a*u + b

    u_ = u[n-1]
    k = 0
    while abs(F(u_)) > eps_r and k < max_iter:
        u_ = u_ - F(u_)/dF(u_)
        k += 1
    u[n] = u_
    iterations.append(k)
return u, iterations

```

The Crank-Nicolson method utilizing a linearization based on the geometric mean gives a simpler algorithm:

```

def CN_logistic(u0, dt, Nt):
    u = np.zeros(Nt+1)
    u[0] = u0
    for n in range(0, Nt):
        u[n+1] = (1 + 0.5*dt)/(1 + dt*u[n] - 0.5*dt)*u[n]
    return u

```

We may run experiments with the model problem (1) and the different strategies for dealing with nonlinearities as described above. For a quite coarse time resolution,  $\Delta t = 0.9$ , use of a tolerance  $\epsilon_r = 0.1$  in the stopping criterion introduces an iteration error, especially in the Picard iterations, that is visibly much larger than the time discretization error due to a large  $\Delta t$ . This is illustrated by comparing the upper two plots in Figure 1. The one to the right as a stricter tolerance  $\epsilon = 10^{-3}$ , which leads to all the curves corresponding to Picard and Newton iteration to be on top of each other (and no changes can be usually observed by reducing  $\epsilon_r$  further). The reason why Newton's method does much better than Picard iteration in the upper left plot is that Newton's method with one step comes far below the  $\epsilon_r$  tolerance, while the Picard iteration needs on average 7 iterations to bring the residual down to  $\epsilon_r = 10^{-1}$ , which gives insufficient accuracy in the solution of the nonlinear equation. It is obvious that the Picard1 method gives significant errors in addition to the time discretization unless the time step is as small as in the lower right plot.

The *BE exact* curve corresponds to using the exact solution of the quadratic equation at each time level, so this curve is only affected by the Backward Euler time discretization. The *CN gm* curve corresponds to the theoretically more accurate Crank-Nicolson discretization, combined with a geometric mean for linearization. This curve appears as more accurate, especially if we take the plot in the lower right with a small  $\Delta t$  and an appropriately small  $\epsilon_r$  value as the exact curve.

When it comes to the need for iterations, Figure 2 displays the number of iterations required at each time level for Newton's method and Picard iteration.

The smaller  $\Delta t$  is, the better starting value we have for the iteration, faster the convergence is. With  $\Delta t = 0.9$  Picard iteration requires on average 32 iterations per time step, but this number is dramatically reduced if  $\omega$  is reduced.

However, introducing relaxation and a parameter  $\omega = 0.8$  immediately reduces the average of 32 to 7, indicating that for the large  $\Delta t = 0.9$ , iteration takes too long steps. An approximately optimal value for  $\omega$  in this case is 0.5, which results in an average of only 2 iterations! Even more dramatic impact of  $\omega$  appears when  $\Delta t = 1$ : Picard iteration does not converge at all, but  $\omega = 0.5$  again brings the average number of iterations down to 2.

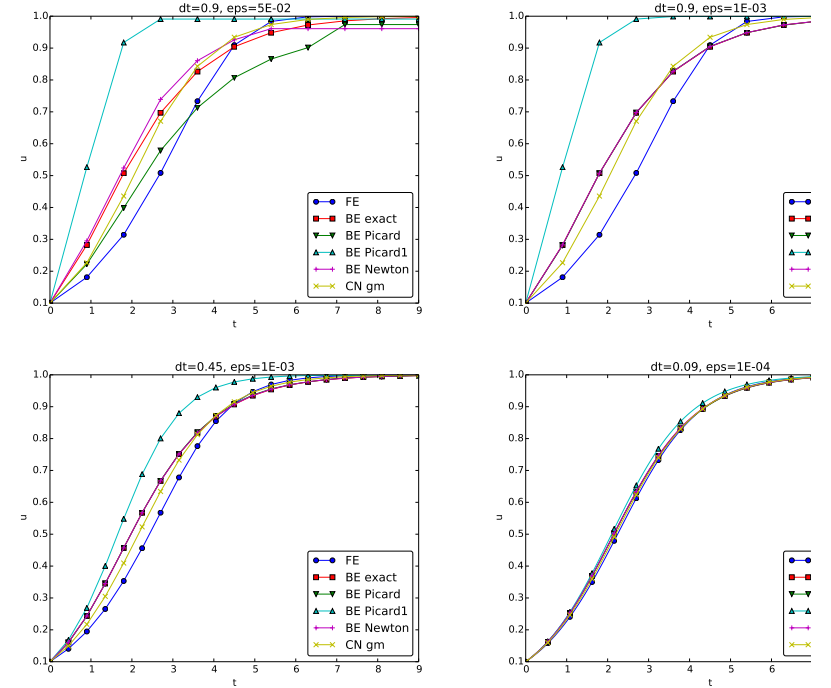


Figure 1: The impact of solution strategies and for four different time step lengths on the solution.

Experiments with this program reveal the relative performance of the methods as summarized in the table below. The Picard and Newton columns refer to the typical number of iterations with these methods before the curve starts to diverge and the number of iterations is significantly reduced since the solution of the nonlinear algebraic equation is very close to the starting value for the iteration (the solution at the previous time level). Increasing  $\Delta t$  moves the starting value further away from the solution of the nonlinear equation and one expects an increase in the number of iterations. Picard iteration is very much more sensitive to  $\Delta t$  than Newton's method.

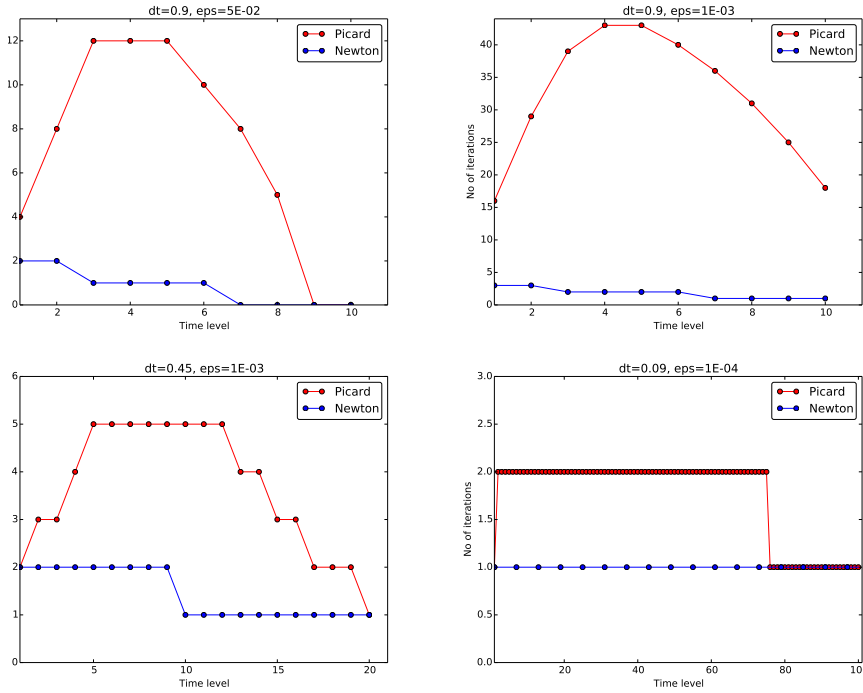


Figure 2: Comparison of the number of iterations at various time levels for Picard and Newton iteration.

the size of  $\Delta t$  than Newton's method. The tolerance  $\epsilon_r$  in residual-based stopping criterion takes on a low and high value in the experiments.

$\Delta t$	$\epsilon_r$	Picard	Newton
0.2	$10^{-7}$	5	2
0.2	$10^{-3}$	2	1
0.4	$10^{-7}$	12	3
0.4	$10^{-3}$	4	2
0.8	$10^{-7}$	58	3
0.8	$10^{-3}$	4	2

**Remark.** The simple Crank-Nicolson method with a geometric mean for the quadratic nonlinearity gives visually more accurate solutions than the Backward Euler discretization. Even with a tolerance of  $\epsilon_r = 10^{-3}$ , all the methods for treating the nonlinearities in the Backward Euler discretization gives graphs that cannot be distinguished. So for accuracy in this problem, the time discretization is much more crucial than  $\epsilon_r$ . Ideally, one should estimate the error in the time discretization, as the solution progresses, and set  $\epsilon_r$  accordingly.

## 1.9 Generalization to a general nonlinear ODE

Let us see how the various methods in the previous sections can be applied to the more generic model

$$u' = f(u, t),$$

where  $f$  is a nonlinear function of  $u$ .

**Explicit time discretization.** Explicit ODE methods like the Forward scheme, Runge-Kutta methods, Adams-Bashforth methods all evaluate time levels where  $u$  is already computed, so nonlinearities in  $f$  do not present difficulties.

**Backward Euler discretization.** Approximating  $u'$  by a backward difference leads to a Backward Euler scheme, which can be written as

$$F(u^n) = u^n - \Delta t f(u^n, t_n) - u^{n-1} = 0,$$

or alternatively

$$F(u) = u - \Delta t f(u, t_n) - u^{(1)} = 0.$$

A simple Picard iteration, not knowing anything about the nonlinear structure of  $f$ , must approximate  $f(u, t_n)$  by  $f(u^-, t_n)$ :

$$\hat{F}(u) = u - \Delta t f(u^-, t_n) - u^{(1)}.$$

The iteration starts with  $u^- = u^{(1)}$  and proceeds with repeating

$$u^* = \Delta t f(u^-, t_n) + u^{(1)}, \quad u = \omega u^* + (1 - \omega)u^-, \quad u^- \leftarrow u,$$

until a stopping criterion is fulfilled.

Newton's method requires the computation of the derivative

$$F'(u) = 1 - \Delta t \frac{\partial f}{\partial u}(u, t_n).$$

Starting with the solution at the previous time level,  $u^- = u^{(1)}$ , we can use the standard formula

$$u = u^- - \omega \frac{F(u^-)}{F'(u^-)} = u^- - \omega \frac{u^{(1)} + \Delta t f(u^-, t_n)}{1 - \Delta t \frac{\partial f}{\partial u}(u^-, t_n)}.$$

The geometric mean trick cannot be used unless we know that  $f$  has a specific structure with quadratic expressions in  $u$ .

**Crank-Nicolson discretization.** The standard Crank-Nicolson scheme with arithmetic mean approximation of  $f$  takes the form

$$\frac{u^{n+1} - u^n}{\Delta t} = \frac{1}{2}(f(u^{n+1}, t_{n+1}) + f(u^n, t_n)).$$

We can write the scheme as a nonlinear algebraic equation

$$F(u) = u - u^{(1)} - \Delta t \frac{1}{2} f(u, t_{n+1}) - \Delta t \frac{1}{2} f(u^{(1)}, t_n) = 0. \quad (12)$$

Picard iteration scheme must in general employ the linearization,

$$\hat{F}(u) = u - u^{(1)} - \Delta t \frac{1}{2} f(u^-, t_{n+1}) - \Delta t \frac{1}{2} f(u^{(1)}, t_n),$$

while Newton's method can apply the general formula (11) with  $F(u)$  given in (12) and

$$F'(u) = 1 - \frac{1}{2} \Delta t \frac{\partial f}{\partial u}(u, t_{n+1}).$$

## 5 Systems of nonlinear algebraic equations

Implicit time discretization methods for a system of ODEs, or a PDE, lead to *systems* of nonlinear algebraic equations, written compactly as

$$F(u) = 0,$$

where  $u$  is a vector of unknowns  $u = (u_0, \dots, u_N)$ , and  $F$  is a vector function:  $F = (F_0, \dots, F_N)$ . Sometimes the equation system has a special structure because of the underlying problem, e.g.,

$$A(u)u = b(u),$$

with  $A(u)$  as an  $(N+1) \times (N+1)$  matrix function of  $u$  and  $b$  as a vector function:  $b = (b_0, \dots, b_N)$ .

We shall next explain how Picard iteration and Newton's method can be applied to systems like  $F(u) = 0$  and  $A(u)u = b(u)$ . The exposition has a focus on ideas and practical computations. More theoretical considerations, including quite general results on convergence properties of these methods, can be found in Kelley [1].

### 5.1 Picard iteration

We cannot apply Picard iteration to nonlinear equations unless there is some special structure. For the commonly arising case  $A(u)u = b(u)$  we can linearize the product  $A(u)u$  to  $A(u^-)u$  and  $b(u)$  as  $b(u^-)$ . That is, we use the most previously computed approximation in  $A$  and  $b$  to arrive at a *linear system* for

$$A(u^-)u = b(u^-).$$

A relaxed iteration takes the form

$$A(u^-)u^* = b(u^-), \quad u = \omega u^* + (1 - \omega)u^-.$$

In other words, we solve a system of nonlinear algebraic equations as a set of linear systems.

#### Algorithm for relaxed Picard iteration.

Given  $A(u)u = b(u)$  and an initial guess  $u^-$ , iterate until convergence

1. solve  $A(u^-)u^* = b(u^-)$  with respect to  $u^*$
2.  $u = \omega u^* + (1 - \omega)u^-$
3.  $u^- \leftarrow u$

### 2.2 Newton's method

The natural starting point for Newton's method is the general nonlinear equation  $F(u) = 0$ . As for a scalar equation, the idea is to approximate  $F$  at a known value  $u^-$  by a linear function  $\hat{F}$ , calculated from the first two terms of a Taylor expansion of  $F$ . In the multi-variate case these two terms become

$$F(u^-) + J(u^-) \cdot (u - u^-),$$

where  $J$  is the *Jacobian* of  $F$ , defined by

$$J_{i,j} = \frac{\partial F_i}{\partial u_j}.$$

So, the original nonlinear system is approximated by

$$\hat{F}(u) = F(u^-) + J(u^-) \cdot (u - u^-) = 0,$$

which is linear in  $u$  and can be solved in a two-step procedure: first solve  $J\delta u = -F(u^-)$  with respect to the vector  $\delta u$  and then update  $u = u^- + \delta u$ . The relaxation parameter can easily be incorporated:

$$u = \omega(u^- + \delta u) + (1 - \omega)u^- = u^- + \omega\delta u.$$



**Algorithm for Newton's method.**

Given  $F(u) = 0$  and an initial guess  $u^-$ , iterate until convergence:

1. solve  $J\delta u = -F(u^-)$  with respect to  $\delta u$
2.  $u = u^- + \omega\delta u$
3.  $u^- \leftarrow u$

For the special system with structure  $A(u)u = b(u)$ ,

$$F_i = \sum_k A_{i,k}(u)u_k - b_i(u),$$

one gets

$$J_{i,j} = \sum_k \frac{\partial A_{i,k}}{\partial u_j} u_k + A_{i,j} - \frac{\partial b_i}{\partial u_j}. \quad (13)$$

We realize that the Jacobian needed in Newton's method consists of  $A(u^-)$  as in the Picard iteration plus two additional terms arising from the differentiation. Using the notation  $A'(u)$  for  $\partial A/\partial u$  (a quantity with three indices:  $\partial A_{i,k}/\partial u_j$ ), and  $b'(u)$  for  $\partial b/\partial u$  (a quantity with two indices:  $\partial b_i/\partial u_j$ ), we can write the linear system to be solved as

$$(A + A'u + b')\delta u = -Au + b,$$

or

$$(A(u^-) + A'(u^-)u^- + b'(u^-))\delta u = -A(u^-)u^- + b(u^-).$$

Rearranging the terms demonstrates the difference from the system solved in each Picard iteration:

$$\underbrace{A(u^-)(u^- + \delta u) - b(u^-)}_{\text{Picard system}} + \gamma(A'(u^-)u^- + b'(u^-))\delta u = 0.$$

Here we have inserted a parameter  $\gamma$  such that  $\gamma = 0$  gives the Picard system and  $\gamma = 1$  gives the Newton system. Such a parameter can be handy in software to easily switch between the methods.

**Combined algorithm for Picard and Newton iteration.**

Given  $A(u)$ ,  $b(u)$ , and an initial guess  $u^-$ , iterate until convergence:

1. solve  $(A + \gamma(A'(u^-)u^- + b'(u^-)))\delta u = -A(u^-)u^- + b(u^-)$  with respect to  $\delta u$
2.  $u = u^- + \omega\delta u$
3.  $u^- \leftarrow u$

$\gamma = 1$  gives a Newton method while  $\gamma = 0$  corresponds to Picard iteration.

**2.3 Stopping criteria**

Let  $\|\cdot\|$  be the standard Euclidean vector norm. Four termination criteria are much in use:

- Absolute change in solution:  $\|u - u^-\| \leq \epsilon_u$
- Relative change in solution:  $\|u - u^-\| \leq \epsilon_u \|u_0\|$ , where  $u_0$  denotes the initial start value of  $u^-$  in the iteration
- Absolute residual:  $\|F(u)\| \leq \epsilon_r$
- Relative residual:  $\|F(u)\| \leq \epsilon_r \|F(u_0)\|$

To prevent divergent iterations to run forever, one terminates the iteration when the current number of iterations  $k$  exceeds a maximum value  $k_{\max}$ .

The relative criteria are most used since they are not sensitive to the characteristic size of  $u$ . Nevertheless, the relative criteria can be misleading if the initial start value for the iteration is very close to the solution, since an unnecessary reduction in the error measure is enforced. In such cases the absolute criteria work better. It is common to combine the absolute and relative norm of the size of the residual, as in

$$\|F(u)\| \leq \epsilon_{rr} \|F(u_0)\| + \epsilon_{ra},$$

where  $\epsilon_{rr}$  is the tolerance in the relative criterion and  $\epsilon_{ra}$  is the tolerance in the absolute criterion. With a very good initial guess for the iteration (typical solution of a differential equation at the previous time level), the term  $\|F(u_0)\|$  is small and  $\epsilon_{ra}$  is the dominating tolerance. Otherwise,  $\epsilon_{rr} \|F(u_0)\|$  and the relative criterion dominates.

With the change in solution as criterion we can formulate a combined absolute and relative measure of the change in the solution:

$$\|\delta u\| \leq \epsilon_{ur} \|u_0\| + \epsilon_{ua},$$

The ultimate termination criterion, combining the residual and the change in solution with a test on the maximum number of iterations allowed, is expressed as

$$\|F(u)\| \leq \epsilon_{rr} \|F(u_0)\| + \epsilon_{ra} \quad \text{or} \quad \|\delta u\| \leq \epsilon_{ur} \|u_0\| + \epsilon_{ua} \quad \text{or} \quad k > k_{\max}$$

#### .4 Example: A nonlinear ODE model from epidemiology

he simplest model spreading of a disease, such as a flu, takes the form of a  $\times 2$  ODE system

$$S' = -\beta SI, \quad (17)$$

$$I' = \beta SI - \nu I, \quad (18)$$

here  $S(t)$  is the number of people who can get ill (susceptibles) and  $I(t)$  is the number of people who are ill (infected). The constants  $\beta > 0$  and  $\nu > 0$  must be given along with initial conditions  $S(0)$  and  $I(0)$ .

**mplicit time discretization.** A Crank-Nicolson scheme leads to a  $2 \times 2$  system of nonlinear algebraic equations in the unknowns  $S^{n+1}$  and  $I^{n+1}$ :

$$\frac{S^{n+1} - S^n}{\Delta t} = -\beta[S I]^{n+\frac{1}{2}} \approx -\frac{\beta}{2}(S^n I^n + S^{n+1} I^{n+1}), \quad (19)$$

$$\frac{I^{n+1} - I^n}{\Delta t} = \beta[S I]^{n+\frac{1}{2}} - \nu I^{n+\frac{1}{2}} \approx \frac{\beta}{2}(S^n I^n + S^{n+1} I^{n+1}) - \frac{\nu}{2}(I^n + I^{n+1}). \quad (20)$$

Introducing  $S$  for  $S^{n+1}$ ,  $S^{(1)}$  for  $S^n$ ,  $I$  for  $I^{n+1}$ ,  $I^{(1)}$  for  $I^n$ , we can rewrite the system as

$$F_S(S, I) = S - S^{(1)} + \frac{1}{2}\Delta t\beta(S^{(1)}I^{(1)} + SI) = 0, \quad (21)$$

$$F_I(S, I) = I - I^{(1)} - \frac{1}{2}\Delta t\beta(S^{(1)}I^{(1)} + SI) - \frac{1}{2}\Delta t\nu(I^{(1)} + I) = 0. \quad (22)$$

**Picard iteration.** We assume that we have approximations  $S_-$  and  $I_-$  to  $S$  and  $I$ . A way of linearizing the only nonlinear term  $SI$  is to write  $I_-S$  in the  $F_S = 0$  equation and  $S_-I$  in the  $F_I = 0$  equation, which also *decouples* the equations. Solving the resulting linear equations with respect to the unknowns  $S$  and  $I$  gives

$$S = \frac{S^{(1)} - \frac{1}{2}\Delta t\beta S^{(1)}I^{(1)}}{1 + \frac{1}{2}\Delta t\beta I_-},$$

$$I = \frac{I^{(1)} + \frac{1}{2}\Delta t\beta S^{(1)}I^{(1)}}{1 - \frac{1}{2}\Delta t\beta S_- + \nu}.$$

Before a new iteration, we must update  $S^- \leftarrow S$  and  $I^- \leftarrow I$ .

**Newton's method.** The nonlinear system (21)-(22) can be written as  $F(u) = 0$  with  $F = (F_S, F_I)$  and  $u = (S, I)$ . The Jacobian becomes

$$J = \begin{pmatrix} \frac{\partial}{\partial S} F_S & \frac{\partial}{\partial I} F_S \\ \frac{\partial}{\partial S} F_I & \frac{\partial}{\partial I} F_I \end{pmatrix} = \begin{pmatrix} 1 + \frac{1}{2}\Delta t\beta I_- & \frac{1}{2}\Delta t\beta \\ -\frac{1}{2}\Delta t\beta S_- & 1 - \frac{1}{2}\Delta t\beta I_- - \frac{1}{2}\Delta t\nu \end{pmatrix}$$

The Newton system  $J(u^-)\delta u = -F(u^-)$  to be solved in each iteration

$$\begin{pmatrix} 1 + \frac{1}{2}\Delta t\beta I_- & \frac{1}{2}\Delta t\beta S_- \\ -\frac{1}{2}\Delta t\beta S_- & 1 - \frac{1}{2}\Delta t\beta I_- - \frac{1}{2}\Delta t\nu \end{pmatrix} \begin{pmatrix} \delta S \\ \delta I \end{pmatrix} = \begin{pmatrix} S_- - S^{(1)} + \frac{1}{2}\Delta t\beta(S^{(1)}I^{(1)} + S_-I_-) \\ I_- - I^{(1)} - \frac{1}{2}\Delta t\beta(S^{(1)}I^{(1)} + S_-I_-) - \frac{1}{2}\Delta t\nu(I^{(1)} + I_-) \end{pmatrix}$$

**Remark.** For this particular system of ODEs, explicit time integration works very well. The 4-th order Runge-Kutta method is an excellent choice between high accuracy, high efficiency, and simplicity.

### 3 Linearization at the differential equation

The attention is now turned to nonlinear partial differential equations and application of the techniques explained above for ODEs. The model is a nonlinear diffusion equation

$$\begin{aligned} \frac{\partial u}{\partial t} &= \nabla \cdot (\alpha(u)\nabla u) + f(u), & \mathbf{x} \in \Omega, \quad t \in (0, T], \\ -\alpha(u)\frac{\partial u}{\partial n} &= g, & \mathbf{x} \in \partial\Omega_N, \quad t \in (0, T], \\ u &= u_0, & \mathbf{x} \in \partial\Omega_D, \quad t \in (0, T]. \end{aligned}$$

Our aim is to discretize the problem in time and then present techniques for linearizing the time-discrete PDE problem “at the PDE level” so we transform the nonlinear stationary PDE problems at each time level into a sequence of linear PDE problems, which can be solved using any method for linear PDEs. This strategy avoids the solution of nonlinear algebraic equations. In Section 4 we shall take the opposite (and more common) approach: discretize the nonlinear problem in time and space first, and then solve the resulting nonlinear algebraic equations at each time level by the method of Section 2.

## .1 Explicit time integration

The nonlinearities in the PDE are trivial to deal with if we choose an explicit time integration method for (23), such as the Forward Euler method:

$$[D_t^+ u = \nabla \cdot (\alpha(u) \nabla u) + f(u)]^n,$$

can be written out,

$$\frac{u^{n+1} - u^n}{\Delta t} = \nabla \cdot (\alpha(u^n) \nabla u^n) + f(u^n),$$

which is a linear equation in the unknown  $u^{n+1}$  with solution

$$u^{n+1} = u^n + \Delta t \nabla \cdot (\alpha(u^n) \nabla u^n) + \Delta t f(u^n).$$

The disadvantage with this discretization is usually thought to be the stability criterion

$$\Delta t \leq \frac{1}{\max \alpha} (\Delta x^2 + \Delta y^2 + \Delta z^2),$$

or the case  $f = 0$  and a standard 2nd-order finite difference discretization in space with mesh cell sizes  $\Delta x$ ,  $\Delta y$ , and  $\Delta z$  in the various spatial directions.

## .2 Backward Euler scheme and Picard iteration

The Backward Euler scheme for (23) reads

$$[D_t^- u = \nabla \cdot (\alpha(u) \nabla u) + f(u)]^n.$$

can be written out,

$$\frac{u^n - u^{n-1}}{\Delta t} = \nabla \cdot (\alpha(u^n) \nabla u^n) + f(u^n). \quad (26)$$

This is a nonlinear, stationary PDE for the unknown function  $u^n(\mathbf{x})$ . We introduce a Picard iteration with  $k$  as iteration counter. A typical linearization of the  $\nabla \cdot (\alpha(u^n) \nabla u^n)$  term in iteration  $k + 1$  is to use the previously computed  $u^{n,k}$  approximation in the diffusion coefficient:  $\alpha(u^{n,k})$ . The nonlinear source term is treated similarly:  $f(u^{n,k})$ . The unknown function  $u^{n,k+1}$  then fulfills the linear PDE

$$\frac{u^{n,k+1} - u^{n-1}}{\Delta t} = \nabla \cdot (\alpha(u^{n,k}) \nabla u^{n,k+1}) + f(u^{n,k}). \quad (27)$$

The initial guess for the Picard iteration at this time level can be taken as the solution at the previous time level:  $u^{n,0} = u^{n-1}$ .

We can alternatively apply the implementation-friendly notation where  $u$  corresponds to the unknown we want to solve for, i.e.,  $u^{n,k+1}$  above, and  $u^-$  is the most recently computed value,  $u^{n,k}$  above. Moreover,  $u^{(1)}$  denotes the

unknown function at the previous time level,  $u^{n-1}$  above. The PDE to be solved in a Picard iteration then looks like

$$\frac{u - u^{(1)}}{\Delta t} = \nabla \cdot (\alpha(u^-) \nabla u) + f(u^-).$$

At the beginning of the iteration we start with the value from the previous level:  $u^- = u^{(1)}$ , and after each iteration,  $u^-$  is updated to  $u$ .

## 3.3 Backward Euler scheme and Newton's method

At time level  $n$  we have to solve the stationary PDE (26), this time with Newton's method. Normally, Newton's method is defined for systems of *algebraic equations* but the idea of the method can be applied at the PDE level too.

**Linearization via Taylor expansions.** Let  $u^{n,k}$  be an approximation of the unknown  $u^n$ . We seek a better approximation on the form

$$u^n = u^{n,k} + \delta u.$$

The idea is to insert (29) in (26), Taylor expand the nonlinearities and keep the terms that are linear in  $\delta u$ . Then we can solve a linear PDE for the correction  $\delta u$  and use (29) to find a new approximation  $u^{n,k+1} = u^{n,k} + \delta u$ .

Inserting (29) in (26) gives

$$\frac{u^{n,k} + \delta u - u^{n-1}}{\Delta t} = \nabla \cdot (\alpha(u^{n,k} + \delta u) \nabla (u^{n,k} + \delta u)) + f(u^{n,k} + \delta u)$$

We can Taylor expand  $\alpha(u^{n,k} + \delta u)$  and  $f(u^{n,k} + \delta u)$ :

$$\begin{aligned} \alpha(u^{n,k} + \delta u) &= \alpha(u^{n,k}) + \frac{d\alpha}{du}(u^{n,k})\delta u + \mathcal{O}(\delta u^2) \approx \alpha(u^{n,k}) + \alpha'(u^{n,k})\delta u, \\ f(u^{n,k} + \delta u) &= f(u^{n,k}) + \frac{df}{du}(u^{n,k})\delta u + \mathcal{O}(\delta u^2) \approx f(u^{n,k}) + f'(u^{n,k})\delta u, \end{aligned}$$

Inserting the linear approximations of  $\alpha$  and  $f$  in (30) results in

$$\begin{aligned} \frac{u^{n,k} + \delta u - u^{n-1}}{\Delta t} &= \nabla \cdot (\alpha(u^{n,k}) \nabla u^{n,k}) + f(u^{n,k}) + \\ &\quad \nabla \cdot (\alpha'(u^{n,k}) \delta u \nabla u^{n,k}) + \\ &\quad \nabla \cdot (\alpha'(u^{n,k}) \delta u \nabla \delta u) + f'(u^{n,k}) \delta u \end{aligned}$$

The term  $\alpha'(u^{n,k}) \delta u \nabla \delta u$  is  $\mathcal{O}(\delta u^2)$  and therefore omitted. Reorganizing the equation gives a PDE for  $\delta u$  that we can write in short form as

$$\delta F(\delta u; u^{n,k}) = -F(u^{n,k}),$$

here

$$F(u^{n,k}) = \frac{u^{n,k} - u^{n-1}}{\Delta t} - \nabla \cdot (\alpha(u^{n,k}) \nabla u^{n,k}) + f(u^{n,k}), \quad (32)$$

$$\begin{aligned} \delta F(\delta u; u^{n,k}) = & -\frac{1}{\Delta t} \delta u + \nabla \cdot (\alpha(u^{n,k}) \nabla \delta u) + \\ & \nabla \cdot (\alpha'(u^{n,k}) \delta u \nabla u^{n,k}) + f'(u^{n,k}) \delta u. \end{aligned} \quad (33)$$

Note that  $\delta F$  is a linear function of  $\delta u$ , and  $F$  contains only terms that are known, such that the PDE for  $\delta u$  is indeed linear.

#### Observations.

The notational form  $\delta F = -F$  resembles the Newton system  $J\delta u = -F$  for systems of algebraic equations, with  $\delta F$  as  $J\delta u$ . The unknown vector in a linear system of algebraic equations enters the system as a linear operator in terms of a matrix-vector product ( $J\delta u$ ), while at the PDE level we have a linear differential operator instead ( $\delta F$ ).

**Similarity with Picard iteration.** We can rewrite the PDE for  $\delta u$  in a slightly different way too if we define  $u^{n,k} + \delta u$  as  $u^{n,k+1}$ .

$$\begin{aligned} \frac{u^{n,k+1} - u^{n-1}}{\Delta t} = & \nabla \cdot (\alpha(u^{n,k}) \nabla u^{n,k+1}) + f(u^{n,k}) \\ & + \nabla \cdot (\alpha'(u^{n,k}) \delta u \nabla u^{n,k}) + f'(u^{n,k}) \delta u. \end{aligned} \quad (34)$$

Note that the first line is the same PDE as arise in the Picard iteration, while the remaining terms arise from the differentiations that are an inherent ingredient in Newton's method.

**Implementation.** For coding we want to introduce  $u$  for  $u^n$ ,  $u^-$  for  $u^{n,k}$  and  $u^{(1)}$  for  $u^{n-1}$ . The formulas for  $F$  and  $\delta F$  are then more clearly written as

$$F(u^-) = \frac{u^- - u^{(1)}}{\Delta t} - \nabla \cdot (\alpha(u^-) \nabla u^-) + f(u^-), \quad (35)$$

$$\begin{aligned} \delta F(\delta u; u^-) = & -\frac{1}{\Delta t} \delta u + \nabla \cdot (\alpha(u^-) \nabla \delta u) + \\ & \nabla \cdot (\alpha'(u^-) \delta u \nabla u^-) + f'(u^-) \delta u. \end{aligned} \quad (36)$$

The form that orders the PDE as the Picard iteration terms plus the method's derivative terms becomes

$$\begin{aligned} \frac{u - u^{(1)}}{\Delta t} = & \nabla \cdot (\alpha(u^-) \nabla u) + f(u^-) + \\ & \gamma (\nabla \cdot (\alpha'(u^-) (u - u^-) \nabla u^-) + f'(u^-) (u - u^-)). \end{aligned}$$

The Picard and full Newton versions correspond to  $\gamma = 0$  and  $\gamma = 1$ , resp

### 3.4 Crank-Nicolson discretization

A Crank-Nicolson discretization of (23) applies a centered difference at

$$[D_t u = \nabla \cdot (\alpha(u) \nabla u) + f(u)]^{n+\frac{1}{2}}.$$

Since  $u$  is not known at  $t_{n+\frac{1}{2}}$  we need to express the terms on the right-hand side via unknowns  $u^n$  and  $u^{n+1}$ . The standard technique is to apply an arithmetic average,

$$u^{n+\frac{1}{2}} \approx \frac{1}{2}(u^n + u^{n+1}).$$

However, with nonlinear terms we have many choices of formulating an arithmetic mean:

$$[f(u)]^{n+\frac{1}{2}} \approx f\left(\frac{1}{2}(u^n + u^{n+1})\right) = [f(\bar{u}^t)]^{n+\frac{1}{2}},$$

$$[f(u)]^{n+\frac{1}{2}} \approx \frac{1}{2}(f(u^n) + f(u^{n+1})) = [\bar{f}(u)]^{n+\frac{1}{2}},$$

$$[\alpha(u) \nabla u]^{n+\frac{1}{2}} \approx \alpha\left(\frac{1}{2}(u^n + u^{n+1})\right) \nabla\left(\frac{1}{2}(u^n + u^{n+1})\right) = \alpha(\bar{u}^t) \nabla \bar{u}^t]^{n+\frac{1}{2}},$$

$$[\alpha(u) \nabla u]^{n+\frac{1}{2}} \approx \frac{1}{2}(\alpha(u^n) + \alpha(u^{n+1})) \nabla\left(\frac{1}{2}(u^n + u^{n+1})\right) = [\bar{\alpha}(u)^t \nabla \bar{u}^t]^{n+\frac{1}{2}},$$

$$[\alpha(u) \nabla u]^{n+\frac{1}{2}} \approx \frac{1}{2}(\alpha(u^n) \nabla u^n + \alpha(u^{n+1}) \nabla u^{n+1}) = [\bar{\alpha}(u) \nabla u]^t]^{n+\frac{1}{2}}.$$

## 4 Discretization of stationary nonlinear differential equations

Section 3 presents methods for linearizing time-discrete PDEs directly discretization in space. We can alternatively carry out the discretization of the time-discrete nonlinear PDE problem and get a system of nonlinear algebraic equations, which can be solved by Picard iteration or Newton's method as presented in Section 2. This latter approach will now be described in detail.

We shall work with the 1D problem

$$-(\alpha(u)u')' + au = f(u), \quad x \in (0, L), \quad \alpha(u(0))u'(0) = C, \quad u(L) = D. \quad (43)$$

his problem is of the same nature as those arising from implicit time integration of a nonlinear diffusion PDE as outlined in Section 3.2 (set  $a = 1/\Delta t$  and let  $(u)$  incorporate the nonlinear source term as well as known terms with the  $m$ -dependent unknown function at the previous time level).

## 1 Finite difference discretizations

he nonlinearity in the differential equation (43) poses no more difficulty than a variable coefficient, as in  $(\alpha(x)u')'$ . We can therefore use a standard approach to discretizing the Laplace term with a variable coefficient:

$$[-D_x \alpha D_x u + au = f]_i.$$

Writing this out for a uniform mesh with points  $x_i = i\Delta x$ ,  $i = 0, \dots, N_x$ , leads to

$$-\frac{1}{\Delta x^2} \left( \alpha_{i+\frac{1}{2}}(u_{i+1} - u_i) - \alpha_{i-\frac{1}{2}}(u_i - u_{i-1}) \right) + au_i = f(u_i). \quad (44)$$

his equation is valid at all the mesh points  $i = 0, 1, \dots, N_x - 1$ . At  $i = N_x$  we have the Dirichlet condition  $u_i = D$ . The only difference from the case with  $(\alpha(x)u')'$  and  $f(x)$  is that now  $\alpha$  and  $f$  are functions of  $u$  and not only on  $x$ :  $(\alpha(u)u')'$  and  $f(u(x))$ .

The quantity  $\alpha_{i+\frac{1}{2}}$ , evaluated between two mesh points, needs a comment. Since  $\alpha$  depends on  $u$  and  $u$  is only known at the mesh points, we need to express  $\alpha_{i+\frac{1}{2}}$  in terms of  $u_i$  and  $u_{i+1}$ . For this purpose we use an arithmetic mean, although a harmonic mean is also common in this context if  $\alpha$  features large jumps. There are two choices of arithmetic means:

$$\alpha_{i+\frac{1}{2}} \approx \alpha\left(\frac{1}{2}(u_i + u_{i+1})\right) = [\alpha(\bar{u}^x)]^{i+\frac{1}{2}}, \quad (45)$$

$$\alpha_{i+\frac{1}{2}} \approx \frac{1}{2}(\alpha(u_i) + \alpha(u_{i+1})) = [\overline{\alpha(u)}^x]^{i+\frac{1}{2}} \quad (46)$$

Equation (44) with the latter approximation then looks like

$$-\frac{1}{2\Delta x^2} ((\alpha(u_i) + \alpha(u_{i+1}))(u_{i+1} - u_i) - (\alpha(u_{i-1}) + \alpha(u_i))(u_i - u_{i-1})) + au_i = f(u_i), \quad (47)$$

or written more compactly,

$$[-D_x \bar{\alpha}^x D_x u + au = f]_i.$$

At mesh point  $i = 0$  we have the boundary condition  $\alpha(u)u' = C$ , discretized by

$$[\alpha(u)D_{2x}u = C]_0,$$

meaning

$$\alpha(u_0) \frac{u_1 - u_{-1}}{2\Delta x} = C.$$

The fictitious value  $u_{-1}$  can be eliminated with the aid of (47) for  $i = 0$ . If (47) should be solved with respect to  $u_{i-1}$  and that value (for  $i = 0$ ) should be inserted in (48), but it is algebraically much easier to do it the other way. Alternatively, one can use a ghost cell  $[-\Delta x, 0]$  and update the  $u_{-1}$  in the ghost cell according to (48) after every Picard or Newton iteration; this approach means that we use a known  $u_{-1}$  value in (47) from the previous iteration.

## 4.2 Solution of algebraic equations

**The structure of the equation system.** The nonlinear algebraic equations (47) are of the form  $A(u)u = b(u)$  with

$$\begin{aligned} A_{i,i} &= \frac{1}{2\Delta x^2} (-\alpha(u_{i-1}) + 2\alpha(u_i) - \alpha(u_{i+1})) + a, \\ A_{i,i-1} &= -\frac{1}{2\Delta x^2} (\alpha(u_{i-1}) + \alpha(u_i)), \\ A_{i,i+1} &= -\frac{1}{2\Delta x^2} (\alpha(u_i) + \alpha(u_{i+1})), \\ b_i &= f(u_i). \end{aligned}$$

The matrix  $A(u)$  is tridiagonal:  $A_{i,j} = 0$  for  $j > i+1$  and  $j < i-1$ .

The above expressions are valid for internal mesh points  $1 \leq i \leq N_x$ . For  $i = 0$  we need to express  $u_{i-1} = u_{-1}$  in terms of  $u_1$  using (48):

$$u_{-1} = u_1 - \frac{2\Delta x}{\alpha(u_0)}.$$

This value must be inserted in  $A_{0,0}$ . The expression for  $A_{i,i+1}$  applies for  $i = 0$  and  $A_{i,i-1}$  does not enter the system when  $i = 0$ .

Regarding the last equation, its form depends on whether we include the Dirichlet condition  $u(L) = D$ , meaning  $u_{N_x} = D$ , in the nonlinear algebraic equation system or not. Suppose we choose  $(u_0, u_1, \dots, u_{N_x-1})$  as unknowns; this is later referred to as *equations without Dirichlet conditions*. The last equation corresponds to  $i = N_x - 1$ . It involves the boundary value  $u_{N_x}$ , which is substituted by  $D$ . If the unknown vector includes the boundary value,  $(u_0, u_1, \dots, u_{N_x})$ , this is later referred to as *equations including Dirichlet conditions*, the last equation corresponds to  $i = N_x$  just involves the unknown  $u_{N_x}$ , and the final equation becomes  $u_{N_x} = D$ , corresponding to  $A_{i,i} = 1$  and  $b_i = D$  for  $i = N_x$ .

**Picard iteration.** The obvious Picard iteration scheme is to use previously computed values of  $u_i$  in  $A(u)$  and  $b(u)$ , as described more in detail in Section 2. The system  $F(u) = 0$  is then solved with respect to  $u$ , where  $F = (F_0, F_1, \dots, F_m)$ ,  $u = (u_0, u_1, \dots, u_m)$ , and the  $F_i$  expression is given above. The index  $m$  is  $N_x$  in equations including the Dirichlet condition and  $N_x - 1$  when the Dirichlet condition is excluded.

To write out the mathematical details, we introduce  $u^-$  as the most recent approximation to solution vector  $u$ , and  $u_i^-$  is the  $i$ -th component in  $u$ , which is the most recently computed value of the unknown  $u_i$ . For the case  $N_x = 2$  we set the following system to solve in case we omit the Dirichlet condition from the system:

$$\begin{pmatrix} \frac{1}{2\Delta x^2}(-\alpha(u_1^-) + 2\alpha(u_0^-) - \alpha(u_1^-)) + a & -\frac{1}{2\Delta x^2}(\alpha(u_0^-) + \alpha(u_1^-)) \\ -\frac{1}{2\Delta x^2}(\alpha(u_0^-) + \alpha(u_1^-)) & \frac{1}{2\Delta x^2}(-\alpha(u_0^-) + 2\alpha(u_1^-) - \alpha(u_2)) + a \end{pmatrix} \begin{pmatrix} u_0 \\ u_1 \end{pmatrix} = \begin{pmatrix} 0 \\ 0 \end{pmatrix}$$

here  $u_{-1}$  must be substituted by (49), and  $u_2$  by  $D$ .

The system with the Dirichlet condition becomes

$$\begin{pmatrix} \frac{1}{2\Delta x^2}(-\alpha(u_1^-) + 2\alpha(u_0^-) - \alpha(u_1^-)) + a & 0 \\ -\frac{1}{2\Delta x^2}(\alpha(u_0^-) + \alpha(u_1^-)) & \frac{1}{2\Delta x^2}(-\alpha(u_0^-) + 2\alpha(u_1^-) - \alpha(u_2)) + a \\ -\frac{1}{2\Delta x^2}(\alpha(u_0^-) + \alpha(u_1^-)) & 0 \end{pmatrix} \begin{pmatrix} u_0 \\ u_1 \\ u_2 \end{pmatrix} = \begin{pmatrix} 0 \\ 0 \\ 0 \end{pmatrix}$$

**Newton's method.** The Jacobian must be derived in order to use Newton's method. Here it means that we need to differentiate  $F(u) = A(u)u - b(u)$  with respect to the unknown parameters  $u_0, u_1, \dots, u_m$  ( $m = N_x$  or  $m = N_x - 1$ , depending on whether the Dirichlet condition is included in the nonlinear system  $F(u) = 0$  or not). Nonlinear equation number  $i$  has the structure

$$F_i = A_{i,i-1}(u_{i-1}, u_i)u_{i-1} + A_{i,i}(u_{i-1}, u_i, u_{i+1})u_i + A_{i,i+1}(u_i, u_{i+1})u_{i+1} - b_i(u_i).$$

Computing the Jacobian requires careful differentiation. For example,

$$\begin{aligned} \frac{\partial}{\partial u_i}(A_{i,i}(u_{i-1}, u_i, u_{i+1})u_i) &= \frac{\partial A_{i,i}}{\partial u_i}u_i + A_{i,i} \frac{\partial u_i}{\partial u_i} \\ &= \frac{\partial}{\partial u_i} \left( \frac{1}{2\Delta x^2}(-\alpha(u_{i-1}) + 2\alpha(u_i) - \alpha(u_{i+1})) + a \right) u_i + \\ &\quad \frac{1}{2\Delta x^2}(-\alpha(u_{i-1}) + 2\alpha(u_i) - \alpha(u_{i+1})) + a \\ &= \frac{1}{2\Delta x^2}(2\alpha'(u_i)u_i - \alpha(u_{i-1}) + 2\alpha(u_i) - \alpha(u_{i+1})) + a. \end{aligned}$$

The complete Jacobian becomes

$$\begin{aligned} J_{i,i} &= \frac{\partial F_i}{\partial u_i} = \frac{\partial A_{i,i-1}}{\partial u_i}u_{i-1} + \frac{\partial A_{i,i}}{\partial u_i}u_i + A_{i,i} + \frac{\partial A_{i,i+1}}{\partial u_i}u_{i+1} - \frac{\partial b_i}{\partial u_i} \\ &= \frac{1}{2\Delta x^2}(-\alpha'(u_i)u_{i-1} + 2\alpha'(u_i)u_i - \alpha(u_{i-1}) + 2\alpha(u_i) - \alpha(u_{i+1}) \\ &\quad a - \frac{1}{2\Delta x^2}\alpha'(u_i)u_{i+1} - b'(u_i)), \\ J_{i,i-1} &= \frac{\partial F_i}{\partial u_{i-1}} = \frac{\partial A_{i,i-1}}{\partial u_{i-1}}u_{i-1} + A_{i-1,i} + \frac{\partial A_{i,i}}{\partial u_{i-1}}u_i - \frac{\partial b_i}{\partial u_{i-1}} \\ &= \frac{1}{2\Delta x^2}(-\alpha'(u_{i-1})u_{i-1} - (\alpha(u_{i-1}) + \alpha(u_i)) + \alpha'(u_{i-1})u_i), \\ J_{i,i+1} &= \frac{\partial A_{i,i+1}}{\partial u_{i-1}}u_{i+1} + A_{i+1,i} + \frac{\partial A_{i,i}}{\partial u_{i+1}}u_i - \frac{\partial b_i}{\partial u_{i+1}} \\ &= \frac{1}{2\Delta x^2}(-\alpha'(u_{i+1})u_{i+1} - (\alpha(u_i) + \alpha(u_{i+1})) + \alpha'(u_{i+1})u_i). \end{aligned}$$

The explicit expression for nonlinear equation number  $i$ ,  $F_i(u_0, u_1, \dots)$ , from moving the  $(u_i)$  term in (47) to the left-hand side:

$$F_i = -\frac{1}{2\Delta x^2}((\alpha(u_i) + \alpha(u_{i+1}))(u_{i+1} - u_i) + au_i - f(u_i)) = 0.$$

At the boundary point  $i = 0$ ,  $u_{-1}$  must be replaced using the form (49). When the Dirichlet condition at  $i = N_x$  is not a part of the equation, the last equation  $F_m = 0$  for  $m = N_x - 1$  involves the quantity  $u_{N_x}$  which must be replaced by  $D$ . If  $u_{N_x}$  is treated as an unknown in the system, equation  $F_m = 0$  has  $m = N_x$  and reads

$$F_{N_x}(u_0, \dots, u_{N_x}) = u_{N_x} - D = 0.$$

Similar replacement of  $u_{-1}$  and  $u_{N_x}$  must be done in the Jacobian for the first and last row. When  $u_{N_x}$  is included as an unknown, the last row in the Jacobian must help implement the condition  $\delta u_{N_x} = 0$ , since we assume that  $u_{N_x}$  is the right Dirichlet value at the beginning of the iteration ( $u_{N_x} = D$ ), a Newton update should be zero for  $i = 0$ , i.e.,  $\delta u_{N_x} = 0$ . This also forces the right-hand side to be  $b_i = 0$ ,  $i = N_x$ .

We have seen, and can see from the present example, that the linear system in Newton's method contains all the terms present in the system that is solved in the Picard iteration method. The extra terms in Newton's method are multiplied by a factor such that it is easy to program one linear system with a factor of 0 or 1 to generate the Picard or Newton system.

### 4.3 Galerkin-type discretizations

For the finite element discretization we first need to derive the variational problem. Let  $V$  be an appropriate function space with basis functions  $\{$

Because of the Dirichlet condition at  $x = L$  we require  $\psi_i(L) = 0$ ,  $i \in \mathcal{I}_s$ . The approximate solution is written as  $u = D + \sum_{j \in \mathcal{I}_s} c_j \psi_j$ , where the term  $D$  can be viewed as a boundary function needed to implement the Dirichlet condition  $u(L) = D$ .

Using Galerkin's method, we multiply the differential equation by any  $v \in V$  and integrate terms with second-order derivatives by parts:

$$\int_0^L \alpha(u) u' v' dx + \int_0^L a u v dx = \int_0^L f(u) v dx + [\alpha(u) u' v]_0^L, \quad \forall v \in V.$$

The Neumann condition at the boundary  $x = 0$  is inserted in the boundary term:

$$[\alpha(u) u' v]_0^L = \alpha(u(L)) u'(L) v(L) - \alpha(u(0)) u'(0) v(0) = 0 - C v(0) = -C v(0).$$

Recall that since  $\psi_i(L) = 0$ , any linear combination  $v$  of the basis functions also vanishes at  $x = L$ :  $v(L) = 0$ .) The variational problem is then: find  $u \in V$  such that

$$\int_0^L \alpha(u) u' v' dx + \int_0^L a u v dx = \int_0^L f(u) v dx - C v(0), \quad \forall v \in V. \quad (51)$$

To derive the algebraic equations, we note that  $\forall v \in V$  is equivalent with  $v = \psi_i$  for  $i \in \mathcal{I}_s$ . Setting  $u = D + \sum_j c_j \psi_j$  and sorting terms results in the linear system

$$\sum_{i \in \mathcal{I}_s} \left( \int_0^L \alpha(D + \sum_{k \in \mathcal{I}_s} c_k \psi_k) \psi_j' \psi_i' dx \right) c_j = \int_0^L f(D + \sum_{k \in \mathcal{I}_s} c_k \psi_k) \psi_i dx - C \psi_i(0), \quad i \in \mathcal{I}_s \quad (52)$$

**Fundamental integration problem.** Methods that use the Galerkin or residual principle face a fundamental difficulty in nonlinear problems: how can we integrate terms like  $\int_0^L \alpha(\sum_k c_k \psi_k) \psi_i' \psi_j' dx$  and  $\int_0^L f(\sum_k c_k \psi_k) \psi_i dx$  when we do not know the  $c_k$  coefficients in the argument of the  $\alpha$  function? We can resort to numerical integration, provided an approximate  $\sum_k c_k \psi_k$  can be used for the argument  $u$  in  $f$  and  $\alpha$ . If we want to derive the structure of the nonlinear algebraic equations, we need to apply numerical integration based on the nodes only and/or the group finite element method.

## 4 Finite element basis functions

Introduction of finite element basis functions  $\varphi_i$  means setting

$$\psi_i = \varphi_{\nu(i)}, \quad i \in \mathcal{I}_s,$$

where degree of freedom number  $\nu(i)$  in the mesh corresponds to node number  $i$  ( $c_i$ ). The expansion of  $u$  can still be

$$u = D + \sum_{j \in \mathcal{I}_s} c_j \varphi_{\nu(j)},$$

but is more common in a finite element context to use a boundary function  $B = \sum_{j \in I_b} U_j \varphi_j$ , where  $U_j$  are prescribed Dirichlet conditions for degree of freedom number  $j$  and  $U_j$  is the corresponding value. In the present context this means

$$u = D \varphi_0 + \sum_{j \in \mathcal{I}_s} c_j \varphi_{j+1}, \quad \mathcal{I}_s = \{0, \dots, N_n - 2\}.$$

In the general case with  $u$  prescribed as  $U_j$  at some nodes  $j \in I_b$ , we set

$$u = \sum_{j \in I_b} U_j \varphi_j + \sum_{j \in \mathcal{I}_s} c_j \varphi_{\nu(j)},$$

where  $c_j = u(x^{\nu(j)})$ . That is,  $\nu(j)$  maps unknown number  $j$  to the corresponding node number  $\nu(j)$  such that  $c_j = u(x^{\nu(j)})$ .

## 4.5 The group finite element method

**Finite element approximation of functions of  $u$ .** Since we already have  $u$  as  $\sum_j \varphi_j u(x_j)$ , we may use the same approximation for other functions. For example,

$$f(u) \approx \sum_j f(x_j) \varphi_j,$$

where  $f(x_j)$  is the value of  $f$  at node  $j$ . Since  $f$  is a function of  $u$ ,  $f(u(x_j))$ . Introducing  $u_j$  as a short form for  $u(x_j)$ , we can write

$$f(u) \approx \sum_j f(u_j) \varphi_j.$$

This approximation is known as the *group finite element method* or the *approximation* technique. The index  $j$  runs over all node numbers in the mesh.

The principal advantages of the group finite element method are two:

1. Complicated nonlinear expressions can be simplified to increase efficiency of numerical computations.
2. One can derive *symbolic form* of the difference equations arising from the finite element method in nonlinear problems. The symbolic form is useful for comparing finite element and finite difference equations of nonlinear differential equation problems.

elow, we shall explore point 2 to see exactly how the finite element method creates more complex expressions in the resulting linear system (the difference equations) than the finite difference method does. It turns out that is very difficult to see what kind of turns in the difference equations that arise from  $\int f(u)\varphi_i dx$  without using the group finite element method or numerical integration utilizing the nodes only.

Note, however, that an expression like  $\int f(u)\varphi_i dx$  causes no problems in a computer program as the integral is calculated by numerical integration using an existing approximation of  $u$  in  $f(u)$  such that the integrand can be sampled at any spatial point.

**simplified problem.** Our aim now is to derive symbolic expressions for the difference equations arising from the finite element method in nonlinear problems and compare the expressions with those arising in the finite difference method. To this, let us simplify the model problem and set  $a = 0$ ,  $\alpha = 1$ ,  $f(u) = u^2$ , and have Neumann conditions at both ends such that we get a very simple nonlinear problem  $-u'' = u^2$ ,  $u'(0) = 1$ ,  $u'(L) = 0$ . The variational form is then

$$\int_0^L u'v' dx = \int_0^L u^2v dx - v(0), \quad \forall v \in V.$$

The term with  $u'v'$  is well known so the only new feature is the term  $\int u^2v dx$ .

To make the distance from finite element equations to finite difference equations as short as possible, we shall substitute  $c_j$  in the sum  $u = \sum_j c_j \varphi_j$  by  $c_j = u(x_j)$  since  $c_j$  is the value of  $u$  at node  $j$ . (In the more general case with Dirichlet conditions as well, we have a sum  $\sum_j c_j \varphi_{\nu(j)}$  where  $c_j$  is replaced by  $(x_{\nu(j)})$ . We can then introduce some other counter  $k$  such that it is meaningful to write  $u = \sum_k u_k \varphi_k$ , where  $k$  runs over appropriate node numbers.) The quantity  $u_j$  in  $\sum_j u_j \varphi_j$  is the same as  $u$  at mesh point number  $j$  in the finite difference method, which is commonly denoted  $u_j$ .

**integrating nonlinear functions.** Consider the term  $\int u^2v dx$  in the variational formulation with  $v = \varphi_i$  and  $u = \sum_k \varphi_k u_k$ :

$$\int_0^L \left( \sum_k u_k \varphi_k \right)^2 \varphi_i dx.$$

Evaluating this integral for P1 elements (see Problem 11) results in the expression

$$\frac{h}{12} (u_{i-1}^2 + 2u_i(u_{i-1} + u_{i+1}) + 6u_i^2 + u_{i+1}^2),$$

to be compared with the simple value  $u_i^2$  that would arise in a finite difference discretization when  $u^2$  is sampled at mesh point  $x_i$ . More complicated  $f(u)$  functions give rise to much more lengthy expressions, if it is possible to carry out the integral symbolically at all.

**Application of the group finite element method.** Let us use the group finite element method to derive the terms in the difference equation corresponding to  $f(u)$  in the differential equation. We have

$$\int_0^L f(u)\varphi_i dx \approx \int_0^L \left( \sum_j \varphi_j f(u_j) \right) \varphi_i dx = \sum_j \left( \int_0^L \varphi_i \varphi_j dx \right) f(u_j)$$

We recognize this expression as the mass matrix  $M$ , arising from  $\int \varphi_i \varphi_j$  times the vector  $f = (f(u_0), f(u_1), \dots)$ :  $Mf$ . The associated terms in the difference equations are, for P1 elements,

$$\frac{h}{6} (f(u_{i-1}) + 4f(u_i) + f(u_{i+1})).$$

Occasionally, we want to interpret this expression in terms of finite differences and to this end a rewrite of this expression is convenient:

$$\frac{h}{6} (f(u_{i-1}) + 4f(u_i) + f(u_{i+1})) = h[f(u) - \frac{h^2}{6} D_x D_x f(u)]_i.$$

That is, the finite element treatment of  $f(u)$  (when using a group finite element method) gives the same term as in a finite difference approach,  $f(u_i)$ , plus a diffusion term which is the 2nd-order discretization of  $\frac{1}{6}h^2 f''(x_i)$ .

We may lump the mass matrix through integration with the Trapezoidal rule so that  $M$  becomes diagonal in the finite element method. In that case the term in the differential equation gives rise to a single term  $hf(u_i)$ , just as in the finite difference method.

## 4.6 Numerical integration of nonlinear terms

Let us reconsider a term  $\int f(u)v dx$  as treated in the previous section. Now we want to integrate this term numerically. Such an approach can be written in easy-to-interpret formulas if we apply a numerical integration rule that samples the integrand at the node points  $x_i$  only, because at such points,  $\varphi_j(x_i) = \delta_{ij}$ , which leads to great simplifications.

The term in question takes the form

$$\int_0^L f\left(\sum_k u_k \varphi_k\right) \varphi_i dx.$$

Evaluation of the integrand at a node  $x_\ell$  leads to a collapse of the sum  $\sum_k$  to one term because

$$\begin{aligned} \sum_k u_k \varphi_k(x_\ell) &= u_\ell. \\ f\left(\sum_k u_k \underbrace{\varphi_k(x_\ell)}_{\delta_{k\ell}}\right) \underbrace{\varphi_i(x_\ell)}_{\delta_{i\ell}} &= f(u_\ell) \delta_{i\ell}, \end{aligned}$$



here we have used the Kronecker delta:  $\delta_{ij} = 0$  if  $i \neq j$  and  $\delta_{ij} = 1$  if  $i = j$ .

Considering the Trapezoidal rule for integration, where the integration points are the nodes, we have

$$\int_0^L f\left(\sum_k u_k \varphi_k\right)(x) \varphi_i(x) dx \approx h \sum_{\ell=0}^{N_n} f(u_\ell) \delta_{i\ell} - \mathcal{C} = hf(u_i).$$

This is the same representation of the  $f$  term as in the finite difference method. The term  $\mathcal{C}$  contains the evaluations of the integrand at the ends with weight  $\frac{1}{2}$ , needed to make a true Trapezoidal rule:

$$\mathcal{C} = \frac{h}{2} f(u_0) \varphi_i(0) + \frac{h}{2} f(u_{N_n-1}) \varphi_i(L).$$

The answer  $hf(u_i)$  must therefore be multiplied by  $\frac{1}{2}$  if  $i = 0$  or  $i = N_n - 1$ . Note that  $\mathcal{C} = 0$  for  $i = 1, \dots, N_n - 2$ .

One can alternatively use the Trapezoidal rule on the reference cell and assemble the contributions. It is a bit more labor in this context, but working on the reference cell is safer as that approach is guaranteed to handle discontinuous derivatives of finite element functions correctly (not important in this particular example), while the rule above was derived with the assumption that  $f$  is continuous at the integration points.

The conclusion is that it suffices to use the Trapezoidal rule if one wants to derive the difference equations in the finite element method and make them similar to those arising in the finite difference method. The Trapezoidal rule has sufficient accuracy for P1 elements, but for P2 elements one should turn to Simpson's rule.

## 7 Finite element discretization of a variable coefficient Laplace term

Turning back to the model problem (43), it remains to calculate the contribution of the  $(\alpha u')'$  and boundary terms to the difference equations. The integral in the variational form corresponding to  $(\alpha u')'$  is

$$\int_0^L \alpha \left( \sum_k c_k \psi_k \right) \psi'_i \psi'_j dx.$$

Numerical integration utilizing a value of  $\sum_k c_k \psi_k$  from a previous iteration can in general be used to compute the integral. Now our aim is to integrate symbolically, as much as we can, to obtain some insight into how the finite element method approximates this term. To be able to derive symbolic expressions, we must either turn to the group finite element method or numerical integration in the node points. Finite element basis functions  $\varphi_i$  are now used.

**Group finite element method.** We set  $\alpha(u) \approx \sum_k \alpha(u_k) \varphi_k$ , and write

$$\int_0^L \alpha \left( \sum_k c_k \varphi_k \right) \varphi'_i \varphi'_j dx \approx \sum_k \underbrace{\left( \int_0^L \varphi_k \varphi'_i \varphi'_j dx \right)}_{L_{i,j,k}} \alpha(u_k) = \sum_k L_{i,j,k} \alpha(u_k)$$

Further calculations are now easiest to carry out in the reference cell. For  $N$  elements we can compute  $L_{i,j,k}$  for the two  $k$  values that are relevant in the reference cell. Turning to local indices, one gets

$$L_{r,s,t}^{(e)} = \frac{1}{2h} \begin{pmatrix} 1 & -1 \\ -1 & 1 \end{pmatrix}, \quad t = 0, 1,$$

where  $r, s, t = 0, 1$  are indices over local degrees of freedom in the reference cell ( $i = q(e, r)$ ,  $j = q(e, s)$ , and  $k = q(e, t)$ ). The sum  $\sum_k L_{i,j,k} \alpha(u_k)$  at the cell level becomes  $\sum_{t=0}^1 L_{r,s,t}^{(e)} \alpha(\tilde{u}_t)$ , where  $\tilde{u}_t$  is  $u(x_{q(e,t)})$ , i.e., the value of  $u$  at node number  $t$  in cell number  $e$ . The element matrix becomes

$$\frac{1}{2} (\alpha(\tilde{u}_0) + \alpha(\tilde{u}_1)) \frac{1}{h} \begin{pmatrix} 1 & -1 \\ -1 & 1 \end{pmatrix}.$$

As usual, we employ a left-to-right numbering of cells and nodes. Row  $n$  in the global matrix gets contributions from the first row of the element matrix in cell  $i - 1$  and the last row of the element matrix in cell  $i$ . In cell  $i - 1$  the sum  $\alpha(\tilde{u}_0) + \alpha(\tilde{u}_1)$  corresponds to  $\alpha(u_{i-1}) + \alpha(u_i)$ . The same becomes  $\alpha(u_i) + \alpha(u_{i+1})$  in cell number  $i$ . We can with this insight add the contributions to row number  $i$  in the global matrix:

$$\frac{1}{2h} (-(\alpha(u_{i-1}) + \alpha(u_i)), \quad \alpha(u_{i-1}) + 2\alpha(u_i) + \alpha(u_{i+1}), \quad \alpha(u_i) + \alpha(u_{i+1})).$$

Multiplying by the vector of unknowns  $u_i$  results in a formula that can be arranged to

$$-\frac{1}{h} \left( \frac{1}{2} (\alpha(u_i) + \alpha(u_{i+1})) (u_{i+1} - u_i) - \frac{1}{2} (\alpha(u_{i-1}) + \alpha(u_i)) (u_i - u_{i-1}) \right)$$

which is nothing but the standard finite difference discretization of  $-(\alpha u')'$  with an arithmetic mean of  $\alpha(u)$  (and the usual factor  $h$  because of the integration in the finite element method).

**Numerical integration at the nodes.** Instead of using the group finite element method and exact integration we can turn to the Trapezoidal rule for computing  $\int_0^L \alpha \left( \sum_k u_k \varphi_k \right) \varphi'_i \varphi'_j dx$ , again at the cell level since that is more convenient when we deal with discontinuous functions  $\varphi'_i$ :

$$\begin{aligned}
\int_{-1}^1 \alpha \left( \sum_t \tilde{u}_t \tilde{\varphi}_t \right) \tilde{\varphi}_r' \tilde{\varphi}_s' \frac{h}{2} dX &= \int_{-1}^1 \alpha \left( \sum_{t=0}^1 \tilde{u}_t \tilde{\varphi}_t \right) \frac{2}{h} \frac{d\tilde{\varphi}_r}{dX} \frac{2}{h} \frac{d\tilde{\varphi}_s}{dX} \frac{h}{2} dX \\
&= \frac{1}{2h} (-1)^r (-1)^s \int_{-1}^1 \alpha \left( \sum_{t=0}^1 u_t \tilde{\varphi}_t(X) \right) dX \\
&\approx \frac{1}{2h} (-1)^r (-1)^s \alpha \left( \sum_{t=0}^1 \tilde{\varphi}_t(-1) \tilde{u}_t \right) + \alpha \left( \sum_{t=0}^1 \tilde{\varphi}_t(1) \tilde{u}_t \right) = \frac{1}{2l}
\end{aligned} \tag{55}$$

$$\tag{56}$$

The element matrix in (56) is identical to the one in (53), showing that the group finite element method and Trapezoidal integration are equivalent with standard finite discretization of a nonlinear Laplace term  $(\alpha(u)u)'$  using an arithmetic mean for  $\alpha$ :  $[D_x \bar{x} D_x u]_i$ .

#### Remark about integration in the physical $x$ coordinate.

We might comment on integration in the physical coordinate system too. The common Trapezoidal rule in Section 4.6 cannot be used to integrate derivatives like  $\varphi_i'$ , because the formula is derived under the assumption of a continuous integrand. One must instead use the more basic version of the Trapezoidal rule where all the trapezoids are summed up. This is straightforward, but I think it is even more straightforward to apply the Trapezoidal rule on the reference cell and assemble the contributions.

The term  $\int auv dx$  in the variational form is linear and gives these terms in the algebraic equations:

$$\frac{ah}{6} (u_{i-1} + 4u_i + u_{i+1}) = ah \left[ u - \frac{h^2}{6} D_x D_x u \right]_i.$$

The final term in the variational form is the Neumann condition at the boundary:  $Cv(0) = C\varphi_i(0)$ . With a left-to-right numbering only  $i = 0$  will give a contribution  $Cv(0) = C\delta_{i0}$  (since  $\varphi_i(0) \neq 0$  only for  $i = 0$ ).

#### Summary.

For the equation

$$-(\alpha(u)u)' + au = f(u),$$

P1 finite elements results in difference equations where

- the term  $-(\alpha(u)u)'$  becomes  $-h[D_x \overline{\alpha(u)} D_x u]_i$  if the group finite element method or Trapezoidal integration is applied,
- $f(u)$  becomes  $hf(u_i)$  with Trapezoidal integration or the “mass matrix” representation  $h[f(u) - \frac{h}{6} D_x D_x f(u)]_i$  if computed by a group finite element method,
- $au$  leads to the “mass matrix” form  $ah[u - \frac{h}{6} D_x D_x u]_i$ .

As we now have explicit expressions for the nonlinear difference equations also in the finite element method, a Picard or Newton method can be derived for the finite difference method. However, our efforts in deriving such forms of the difference equations in the finite element method was motivated by a desire to see how nonlinear terms in differential equations make the finite element and difference method different. For practical calculations in computer programs we apply numerical integration, normally the more accurate Gauss-Legendre quadrature rules, to the integrals directly. This allows us to easily evaluate nonlinear algebraic equations for a given numerical approximation of  $u$  (denoted  $u^-$ ). To solve the nonlinear algebraic equations we need to apply the Picard iteration method or Newton’s method to the variational form derived above.

## 4.8 Picard iteration defined from the variational form

We address again the problem (43) with the corresponding variational form. Our aim is to define a Picard iteration based on this variational form with the idea in Picard iteration is to use a previously computed  $u$  value in the nonlinear functions  $\alpha(u)$  and  $f(u)$ . Let  $u^-$  be the available approximation to  $u$  from the previous iteration. The linearized variational form for Picard iteration is

$$\int_0^L (\alpha(u^-)u'v' + auv) dx = \int_0^L f(u^-)v dx - Cv(0), \quad \forall v \in V.$$

This is a linear problem  $a(u, v) = L(v)$  with bilinear and linear forms

$$a(u, v) = \int_0^L (\alpha(u^-)u'v' + auv) dx, \quad L(v) = \int_0^L f(u^-)v dx - Cv(0)$$

Make sure to distinguish the coefficient  $a$  in  $auv$  from the differential coefficient  $a$  in the abstract bilinear form notation  $a(\cdot, \cdot)$ .

The linear system associated with (57) is computed the standard way. Formally, we are back to solving  $-(\alpha(x)u)' + au = f(x)$ . The unknown  $u$  is on the form  $u = B(x) + \sum_{j \in \mathcal{I}_s} c_j \psi_j$ , with  $B(x) = D$  and  $\psi_i = \varphi_{\nu(i)}$ ,  $\nu(i) \in \mathcal{I}_s = \{0, 1, \dots, N = N_n - 2\}$ .

## .9 Newton's method defined from the variational form

Application of Newton's method to the nonlinear variational form (51) arising from the problem (43) requires identification of the nonlinear algebraic equations  $F_i(c_0, \dots, c_N) = 0$ ,  $i \in \mathcal{I}_s$ , and the Jacobian  $J_{i,j} = \partial F_i / \partial c_j$  for  $i, j \in \mathcal{I}_s$ .

The equations  $F_i = 0$  follows from the variational form

$$\int_0^L (\alpha(u)u'v' + auv) dx = \int_0^L f(u)v dx - Cv(0), \quad \forall v \in V,$$

by choosing  $v = \psi_i$ ,  $i \in \mathcal{I}_s$ , and setting  $u = \sum_{j \in \mathcal{I}_s} c_j \psi_j$ , maybe with a boundary condition to incorporate Dirichlet conditions.

With  $v = \psi_i$  we get

$$F_i = \int_0^L (\alpha(u)u'\psi'_i + au\psi_i - f(u)\psi_i) dx + C\psi_i(0) = 0, \quad i \in \mathcal{I}_s. \quad (58)$$

In the differentiations leading to the Jacobian we will frequently use the results

$$\frac{\partial u}{\partial c_j} = \frac{\partial}{\partial c_j} \sum_k c_k \psi_k = \psi_j, \quad \frac{\partial u'}{\partial c_j} = \frac{\partial}{\partial c_j} \sum_k c_k \psi'_k = \psi'_j.$$

The derivation of the Jacobian of (58) goes as

$$\begin{aligned} J_{i,j} &= \frac{\partial F_i}{\partial c_j} = \int_0^L \frac{\partial}{\partial c_j} (\alpha(u)u'\psi'_i + au\psi_i - f(u)\psi_i) dx \\ &= \int_0^L ((\alpha'(u)\frac{\partial u}{\partial c_j}u' + \alpha(u)\frac{\partial u'}{\partial c_j})\psi'_i + a\frac{\partial u}{\partial c_j}\psi_i - f'(u)\frac{\partial u}{\partial c_j}\psi_i) dx \\ &= \int_0^L ((\alpha'(u)\psi_j u' + \alpha(u)\psi'_j \psi'_i + a\psi_j \psi_i - f'(u)\psi_j \psi_i) dx \\ &= \int_0^L (\alpha'(u)u'\psi'_i \psi_j + \alpha(u)\psi'_i \psi'_j + (a - f(u))\psi_i \psi_j) dx \end{aligned} \quad (59)$$

When calculating the right-hand side vector  $F_i$  and the coefficient matrix  $J_{i,j}$  in the linear system to be solved in each Newton iteration, one must use a previously computed  $u$ , denoted by  $u^-$ , for the symbol  $u$  in (58) and (59). With this notation we have

$$F_i = \int_0^L (\alpha(u^-)u'^-\psi'_i + (a - f(u^-))\psi_i) dx - C\psi_i(0), \quad i \in \mathcal{I}_s, \quad (60)$$

$$J_{i,j} = \int_0^L (\alpha'(u^-)u'^-\psi'_i \psi_j + \alpha(u^-)\psi'_i \psi'_j + (a - f(u^-))\psi_i \psi_j) dx, \quad i, j \in \mathcal{I}_s. \quad (61)$$

These expressions can be used for any basis  $\{\psi_i\}_{i \in \mathcal{I}_s}$ . Choosing finite functions for  $\psi_i$ , one will normally want to compute the integral contribution by cell, working in a reference cell. To this end, we restrict the integration cell and transform the cell to  $[-1, 1]$ . The most recently computed approximation  $u^-$  to  $u$  becomes  $\tilde{u}^- = \sum_t \tilde{u}_t^{-1} \tilde{\varphi}_t(X)$  over the reference element, where the value of  $u^-$  at global node (or degree of freedom)  $q(e, t)$  corresponds to local node  $t$  (or degree of freedom). The formulas (60) and (61) then change to

$$\begin{aligned} \tilde{F}_r^{(e)} &= \int_{-1}^1 (\alpha(\tilde{u}^-)\tilde{u}'^-\tilde{\varphi}'_r + (a - f(\tilde{u}^-))\tilde{\varphi}_r) \det J dX - C\tilde{\varphi}_r(0), \\ \tilde{J}_{r,s}^{(e)} &= \int_{-1}^1 (\alpha'(\tilde{u}^-)\tilde{u}'^-\tilde{\varphi}'_r \tilde{\varphi}_s + \alpha(\tilde{u}^-)\tilde{\varphi}'_r \tilde{\varphi}'_s + (a - f(\tilde{u}^-))\tilde{\varphi}_r \tilde{\varphi}_s) \det J dX \end{aligned}$$

with  $r, s \in I_d$  runs over the local degrees of freedom.

Many finite element programs require the user to provide  $F_i$  and  $J_{i,j}$ . Programs, like FEniCS<sup>2</sup>, are capable of automatically deriving  $J_{i,j}$  if specified.

**Dirichlet conditions.** Incorporation of the Dirichlet values by assembling contributions from all degrees of freedom and then modifying the linear system can be obviously applied to Picard iteration as that method involves a linear system. In the Newton system, however, the unknown is a correction  $\delta u$  to the solution. Dirichlet conditions are implemented by inserting the initial guess  $u^-$  for the Newton iteration and implementing  $\delta u_i = 0$  for the known degrees of freedom. The manipulation of the linear system follows the algorithm in the linear problems, the only difference being that the values are zero.

## 5 Multi-dimensional PDE problems

### 5.1 Finite element discretization

The derivation of  $F_i$  and  $J_{i,j}$  in the 1D model problem is easily generalized to multi-dimensional problems. For example, Backward Euler discretization of the PDE

$$u_t = \nabla \cdot (\alpha(u)\nabla u) + f(u),$$

gives the nonlinear time-discrete PDEs

$$u^n - \Delta t \nabla \cdot (\alpha(u^n)\nabla u^n) + f(u^n) = u^{n-1},$$

or with  $u^n$  simply as  $u$  and  $u^{n-1}$  as  $u^{(1)}$ ,

<sup>2</sup><http://fenicsproject.org>

$$u - \Delta t \nabla \cdot (\alpha(u^n) \nabla u) - \Delta t f(u) = u^{(1)}.$$

he variational form, assuming homogeneous Neumann conditions for simplicity, becomes

$$\int_{\Omega} (uv + \Delta t \alpha(u) \nabla u \cdot \nabla v - \Delta t f(u)v - u^{(1)}v) dx. \quad (65)$$

he nonlinear algebraic equations follow from setting  $v = \psi_i$  and using the representation  $u = \sum_k c_k \psi_k$ , which we just write as

$$F_i = \int_{\Omega} (u\psi_i + \Delta t \alpha(u) \nabla u \cdot \nabla \psi_i - \Delta t f(u)\psi_i - u^{(1)}\psi_i) dx. \quad (66)$$

icard iteration needs a linearization where we use the most recent approximation  $u^-$  to  $u$  in  $\alpha$  and  $f$ :

$$F_i \approx \hat{F}_i = \int_{\Omega} (u\psi_i + \Delta t \alpha(u^-) \nabla u \cdot \nabla \psi_i - \Delta t f(u^-)\psi_i - u^{(1)}\psi_i) dx. \quad (67)$$

he equations  $\hat{F}_i = 0$  are now linear and we can easily derive a linear system for the unknown coefficients  $\{c_i\}_{i \in \mathcal{I}_s}$  by inserting  $u = \sum_j c_j \psi_j$ .

In Newton's method we need to evaluate  $F_i$  with the known value  $u^-$  for  $u$ :

$$F_i \approx \hat{F}_i = \int_{\Omega} (u^- \psi_i + \Delta t \alpha(u^-) \nabla u^- \cdot \nabla \psi_i - \Delta t f(u^-)\psi_i - u^{(1)}\psi_i) dx. \quad (68)$$

he Jacobian is obtained by differentiating (66) and using  $\partial u / \partial c_j = \psi_j$ :

$$J_{i,j} = \frac{\partial F_i}{\partial c_j} = \int_{\Omega} (\psi_j \psi_i + \Delta t \alpha'(u) \psi_j \nabla u \cdot \nabla \psi_i + \Delta t \alpha(u) \nabla \psi_j \cdot \nabla \psi_i - \Delta t f'(u) \psi_j \psi_i) dx. \quad (69)$$

he evaluation of  $J_{i,j}$  as the coefficient matrix in the linear system in Newton's method applies the known approximation  $u^-$  for  $u$ :

$$J_{i,j} = \int_{\Omega} (\psi_j \psi_i + \Delta t \alpha'(u^-) \psi_j \nabla u^- \cdot \nabla \psi_i + \Delta t \alpha(u^-) \nabla \psi_j \cdot \nabla \psi_i - \Delta t f'(u^-) \psi_j \psi_i) dx. \quad (70)$$

Hopefully, these examples also show how convenient the notation with  $u$  and  $u^-$  is: the unknown to be computed is always  $u$  and linearization by inserting known (previously computed) values is a matter of adding an underscore. One can take great advantage of this quick notation in software [2].

**Non-homogeneous Neumann conditions.** A natural physical flux condition for the PDE (64) takes the form of a non-homogeneous Neumann condition

$$-\alpha(u) \frac{\partial u}{\partial n} = g, \quad \mathbf{x} \in \partial\Omega_N,$$

where  $g$  is a prescribed function and  $\partial\Omega_N$  is a part of the boundary of the  $\Omega$ . From integrating  $\int_{\Omega} \nabla \cdot (\alpha \nabla u) dx$  by parts, we get a boundary term

$$\int_{\partial\Omega_N} \alpha(u) \frac{\partial u}{\partial n} v ds.$$

Inserting the condition (71) into this term results in an integral over prescribed values:  $-\int_{\partial\Omega_N} g v ds$ . The nonlinearity in the  $\alpha(u)$  coefficient condition therefore does not contribute with a nonlinearity in the variational form.

**Robin conditions.** Heat conduction problems often apply a kind of Newton cooling law, also known as a Robin condition, at the boundary:

$$-\alpha(u) \frac{\partial u}{\partial n} = h_T(u)(u - T_s(t)), \quad \mathbf{x} \in \partial\Omega_R,$$

where  $h_T(u)$  is a heat transfer coefficient between the body ( $\Omega$ ) and the surroundings,  $T_s$  is the temperature of the surroundings, and  $\partial\Omega_R$  is a part of the boundary where this Robin condition applies. The boundary integral (69) becomes

$$\int_{\partial\Omega_R} h_T(u)(u - T_s(t)) v ds,$$

by replacing  $\alpha(u) \partial u / \partial n$  by  $h_T(u - T_s)$ . Often,  $h_T(u)$  can be taken as constant and then the boundary term is linear in  $u$ , otherwise it is nonlinear and contributes to the Jacobian in a Newton method. Linearization in a Picard method typically uses a known value in  $h_T$ , but keeps the  $u$  in  $u - T_s$  as unknown:  $h_T(u^-)(u - T_s(t))$ .

## 5.2 Finite difference discretization

A typical diffusion equation

$$u_t = \nabla \cdot (\alpha(u) \nabla u) + f(u),$$

can be discretized by (e.g.) a Backward Euler scheme, which in 2D can be written

$$[D_t^- u = D_x \bar{\alpha}^x D_x u + D_y \bar{\alpha}^y D_y u + f(u)]_{i,j}^n.$$

We do not dive into details of boundary conditions now. Dirichlet and Neumann conditions are handled as in linear diffusion problems.

Writing the scheme out, putting the unknown values on the left-hand side and known values on the right-hand side, and introducing  $\Delta x = \Delta y = h$  to save some writing, one gets

$$\begin{aligned} u_{i,j}^n &- \frac{\Delta t}{h^2} \left( \frac{1}{2} (\alpha(u_{i,j}^n) + \alpha(u_{i+1,j}^n)) (u_{i+1,j}^n - u_{i,j}^n) \right. \\ &- \frac{1}{2} (\alpha(u_{i-1,j}^n) + \alpha(u_{i,j}^n)) (u_{i,j}^n - u_{i-1,j}^n) \\ &+ \frac{1}{2} (\alpha(u_{i,j}^n) + \alpha(u_{i,j+1}^n)) (u_{i,j+1}^n - u_{i,j}^n) \\ &\left. - \frac{1}{2} (\alpha(u_{i,j-1}^n) + \alpha(u_{i,j}^n)) (u_{i,j}^n - u_{i,j-1}^n) \right) - \Delta t f(u_{i,j}^n) = u_{i,j}^{n-1} \end{aligned}$$

This defines a nonlinear algebraic system  $A(u)u = b(u)$ . A Picard iteration applies old values  $u^-$  in  $\alpha$  and  $f$ , or equivalently, old values for  $u$  in  $A(u)$  and  $b(u)$ . The result is a linear system of the same type as those arising from  $t = \nabla \cdot (\alpha(\mathbf{x}) \nabla u) + f(\mathbf{x}, t)$ .

Newton's method is as usual more involved. We first define the nonlinear algebraic equations to be solved, drop the superscript  $n$ , and introduce  $u^{(1)}$  for  $n=1$ :

$$\begin{aligned} u_{i,j} &- \frac{\Delta t}{h^2} \left( \frac{1}{2} (\alpha(u_{i,j}) + \alpha(u_{i+1,j})) (u_{i+1,j} - u_{i,j}) - \frac{1}{2} (\alpha(u_{i-1,j}) + \alpha(u_{i,j})) (u_{i,j} - u_{i-1,j}) + \right. \\ &\frac{1}{2} (\alpha(u_{i,j}) + \alpha(u_{i,j+1})) (u_{i,j+1} - u_{i,j}) - \frac{1}{2} (\alpha(u_{i,j-1}) + \alpha(u_{i,j})) (u_{i,j} - u_{i,j-1}) \left. \right) \\ &- \Delta t f(u_{i,j}) - u_{i,j}^{n-1} = 0. \end{aligned}$$

It is convenient to work with two indices  $i$  and  $j$  in 2D finite difference discretizations, but it complicates the derivation of the Jacobian, which then gets four indices. The left-hand expression of an equation  $F_{i,j} = 0$  is to be differentiated with respect to each of the unknowns  $u_{r,s}$  (short for  $u_{r,s}^n$ ),  $r \in \mathcal{I}_x$ ,  $s \in \mathcal{I}_y$ ,

$$J_{i,j,r,s} = \frac{\partial F_{i,j}}{\partial u_{r,s}}.$$

Given  $i$  and  $j$ , only a few  $r$  and  $s$  indices give nonzero contribution since  $F_{i,j}$  contains  $u_{i\pm 1,j}$ ,  $u_{i,j\pm 1}$ , and  $u_{i,j}$ . Therefore,  $J_{i,j,r,s}$  is very sparse and we can set  $u$  on the left-hand side of the Newton system as

$$\begin{aligned} J_{i,j,r,s} \delta u_{r,s} &= J_{i,j,i,j} \delta u_{i,j} + J_{i,j,i-1,j} \delta u_{i-1,j} + J_{i,j,i+1,j} \delta u_{i+1,j} + J_{i,j,i,j-1} \delta u_{i,j-1} \\ &+ J_{i,j,i,j+1} \delta u_{i,j+1} \end{aligned}$$

The specific derivatives become

$$\begin{aligned} J_{i,j,i-1,j} &= \frac{\partial F_{i,j}}{\partial u_{i-1,j}} \\ &= \frac{\Delta t}{h^2} (\alpha'(u_{i-1,j}) (u_{i,j} - u_{i-1,j}) + \alpha(u_{i-1,j}) (-1)) \\ J_{i,j,i+1,j} &= \frac{\partial F_{i,j}}{\partial u_{i+1,j}} \\ &= \frac{\Delta t}{h^2} (-\alpha'(u_{i+1,j}) (u_{i+1,j} - u_{i,j}) - \alpha(u_{i+1,j})) \\ J_{i,j,i,j-1} &= \frac{\partial F_{i,j}}{\partial u_{i,j-1}} \\ &= \frac{\Delta t}{h^2} (\alpha'(u_{i,j-1}) (u_{i,j} - u_{i,j-1}) + \alpha(u_{i,j-1}) (-1)) \\ J_{i,j,i,j+1} &= \frac{\partial F_{i,j}}{\partial u_{i,j+1}} \\ &= \frac{\Delta t}{h^2} (-\alpha'(u_{i,j+1}) (u_{i,j+1} - u_{i,j}) - \alpha(u_{i,j+1})) \end{aligned}$$

The  $J_{i,j,i,j}$  entry has a few more terms. Inserting  $u^-$  for  $u$  in the  $J$  form then forming  $J \delta u = -F$  gives the linear system to be solved in each iteration.

### 5.3 Continuation methods

Picard iteration or Newton's method may diverge when solving PDEs with nonlinearities. Relaxation with  $\omega < 1$  may help, but in highly nonlinear problems it can be necessary to introduce a *continuation parameter*  $\Lambda$  in the problem.  $\Lambda = 0$  gives a version of the problem that is easy to solve, while  $\Lambda = 1$  is the problem. The idea is then to increase  $\Lambda$  in steps,  $\Lambda_0 = 0, \Lambda_1 < \dots < \Lambda_n = 1$ , and use the solution from the problem with  $\Lambda_{i-1}$  as initial guess for the iteration with  $\Lambda_i$ .

The continuation method is easiest to understand through an example. Suppose we intend to solve

$$-\nabla \cdot (|\nabla u|^q \nabla u) = f,$$

which is an equation modeling the flow of a non-Newtonian fluid through a channel or pipe. For  $q = 0$  we have the Poisson equation (corresponding to a Newtonian fluid) and the problem is linear. A typical value for pseudoplastic fluids may be  $q_n = -0.8$ . We can introduce the continuation parameter  $\Lambda$  such that  $q = q_n \Lambda$ . Let  $\{\Lambda_\ell\}_{\ell=0}^n$  be the sequence of  $\Lambda$  values in  $[0, 1]$  corresponding  $q$  values  $\{q_\ell\}_{\ell=0}^n$ . We can then solve a sequence of problems

$$-\nabla \cdot (|\nabla u|_\ell^q \nabla u_\ell) = f, \quad \ell = 0, \dots, n,$$

here the initial guess for iterating on  $u^\ell$  is the previously computed solution  $u^{\ell-1}$ . If a particular  $\Lambda_\ell$  leads to convergence problems, one may try a smaller increase in  $\Lambda$ :  $\Lambda_* = \frac{1}{2}(\Lambda_{\ell-1} + \Lambda_\ell)$ , and repeat halving the step in  $\Lambda$  until convergence is reestablished.

## Exercises

### Problem 1: Determine if equations are nonlinear or not

Classify each term in the following equations as linear or nonlinear. Assume that  $a$  and  $b$  are unknown numbers and that  $u$  and  $v$  are unknown functions. All other symbols are known quantities.

1.  $b^2 = 1$
2.  $a + b = 1, a - 2b = 0$
3.  $mu'' + \beta|u'|u' + cu = F(t)$
4.  $u_t = \alpha u_{xx}$
5.  $u_{tt} = c^2 \nabla^2 u$
6.  $u_t = \nabla \cdot (\alpha(u) \nabla u) + f(x, y)$
7.  $u_t + f(u)_x = 0$
8.  $\mathbf{u}_t + \mathbf{u} \cdot \nabla \mathbf{u} = -\nabla p + r \nabla^2 \mathbf{u}, \nabla \cdot \mathbf{u} = 0$  ( $\mathbf{u}$  is a vector field)
9.  $u' = f(u, t)$
10.  $\nabla^2 u = \lambda e^u$

### Exercise 2: Derive a formula

Derive (9) in Section 1.7. Filename: `relaxed_Newton.pdf`.

### Problem 3: Experience the behavior of Newton's method

The program `Newton_demo.py`<sup>3</sup> illustrates graphically each step in Newton's method and is run like

---

```
terminal> python Newton_demo.py f dfdx x0 xmin xmax
```

---

<sup>3</sup>[http://tinyurl.com/nm5587k/nonlin/Newton\\_demo.py](http://tinyurl.com/nm5587k/nonlin/Newton_demo.py)

Use this program to investigate potential problems with Newton's method solving  $e^{-0.5x^2} \cos(\pi x) = 0$ . Try a starting point  $x_0 = 0.8$  and  $x_0 = 0$  and watch the different behavior. Just run

---

```
Terminal> python Newton_demo.py '0.2 + exp(-0.5*x**2)*cos(pi*x)'
'-x*exp(-x**2)*cos(pi*x) - pi*exp(-x**2)*sin(pi*x)' \
0.85 -3 3
```

---

and repeat with 0.85 replaced by 0.8.

### Problem 4: Linearize a nonlinear vibration ODE

Consider a nonlinear vibration problem

$$mu'' + bu'|u'| + s(u) = F(t),$$

where  $m > 0$  is a constant,  $b \geq 0$  is a constant,  $s(u)$  a possibly nonlinear function of  $u$ , and  $F(t)$  is a prescribed function. Such models arise from Newton's law of motion in mechanical vibration problems where  $s(u)$  is a spring or restoring force,  $mu''$  is mass times acceleration, and  $bu'|u'|$  models water or air damping.

Rewrite the equation for  $u$  as a system of two first-order ODEs, and discretize this system by a Crank-Nicolson (centered difference) method. With  $u^n$  we get a nonlinear term  $v^{n+\frac{1}{2}}|v^{n+\frac{1}{2}}|$ . Use both a geometric and an arithmetic average for  $v^{n+\frac{1}{2}}$ . In the latter case, explain how to apply Newton's method to solve the nonlinear equations at each time level.

### Exercise 5: Find the sparsity of the Jacobian

Consider a typical nonlinear Laplace term like  $\nabla \cdot \alpha(u) \nabla u$  discretized by central finite differences. Explain why the Jacobian corresponding to this term has the same sparsity pattern as the matrix associated with the corresponding linear term  $\alpha \nabla^2 u$ .

**Hint.** Set up the unknowns that enter the difference stencil and find the sparsity of the Jacobian that arise from the stencil.

Filename: `nonlin_sparsity_Jacobian.pdf`.

### Exercise 6: Newton's method for linear problems

Suppose we have a linear system  $F(u) = Au - b = 0$ . Apply Newton's method to this system, and show that the method converges in one iteration. Filename: `nonlin_Newton_linear.pdf`.

### Exercise 7: Differentiate a highly nonlinear term

The operator  $\nabla \cdot (\alpha(u) \nabla u)$  with  $\alpha(u) = \|\nabla u\|^q$  appears in several physical problems, especially flow of Non-Newtonian fluids. In a Newton method one has to carry out the differentiation  $\partial \alpha(u) / \partial c_j$ , for  $u = \sum_k c_k \psi_k$ . Show that

$$\frac{\partial}{\partial u_j} \|\nabla u\|^q = q \|\nabla u\|^{q-2} \nabla u \cdot \nabla \psi_j.$$

Filename: `nonlin_differentiate.pdf`.

### Problem 8: Discretize a 1D problem with a nonlinear coefficient

We consider the problem

$$((1 + u^2)u')' = 1, \quad x \in (0, 1), \quad u(0) = u(1) = 0. \quad (75)$$

- a) Discretize (75) by a centered finite difference method on a uniform mesh.
  - b) Discretize (75) by a finite element method with P1 of equal length. Use the Trapezoidal method to compute all integrals. Set up the resulting matrix system.
- Filename: `nonlin_1D_coeff_discretize.pdf`.

### Problem 9: Linearize a 1D problem with a nonlinear coefficient

We have a two-point boundary value problem

$$((1 + u^2)u')' = 1, \quad x \in (0, 1), \quad u(0) = u(1) = 0. \quad (76)$$

- a) Construct a Picard iteration method for (76) without discretizing in space.
  - b) Apply Newton's method to (76) without discretizing in space.
  - c) Discretize (76) by a centered finite difference scheme. Construct a Picard method for the resulting system of nonlinear algebraic equations.
  - d) Discretize (76) by a centered finite difference scheme. Define the system of nonlinear algebraic equations, calculate the Jacobian, and set up Newton's method for solving the system.
- Filename: `nonlin_1D_coeff_linearize.pdf`.

### Problem 10: Finite differences for the 1D Bratu problem

We address the so-called Bratu problem

$$u'' + \lambda e^u = 0, \quad x \in (0, 1), \quad u(0) = u(1) = 0,$$

where  $\lambda$  is a given parameter and  $u$  is a function of  $x$ . This is a widely used problem for studying numerical methods for nonlinear differential equations. The problem (77) has an exact solution

$$u(x) = -2 \ln \left( \frac{\cosh((x - \frac{1}{2})\theta/2)}{\cosh(\theta/4)} \right),$$

where  $\theta$  solves

$$\theta = \sqrt{2\lambda} \cosh(\theta/4).$$

There are two solutions of (77) for  $0 < \lambda < \lambda_c$  and no solution for  $\lambda > \lambda = \lambda_c$  there is one unique solution. The critical value  $\lambda_c$  solves

$$1 = \sqrt{2\lambda_c} \frac{1}{4} \sinh(\theta(\lambda_c)/4).$$

A numerical value is  $\lambda_c = 3.513830719$ .

- a) Discretize (77) by a centered finite difference method.
  - b) Set up the nonlinear equations  $F_i(u_0, u_1, \dots, u_{N_x}) = 0$  from a). Compute the associated Jacobian.
- Filename: `nonlin_1D_Bratu_fd.pdf`.

### Problem 11: Integrate functions of finite element expansions

We shall investigate integrals on the form

$$\int_0^L f\left(\sum_k u_k \varphi_k(x)\right) \varphi_i(x) dx,$$

where  $\varphi_i(x)$  are P1 finite element basis functions and  $u_k$  are unknown coefficients. We introduce more precisely the values of the unknown function  $u$  at nodes  $x_k$ . We introduce node numbering that goes from left to right and also that all cells have the same length  $h$ . Given  $i$ , the integral only gets contributions from  $[x_{i-1}, x_{i+1}]$ . The interval  $\varphi_k(x) = 0$  for  $k < i - 1$  and  $k > i + 1$ , so only three basis functions contribute:

$$\sum_k u_k \varphi_k(x) = u_{i-1} \varphi_{i-1}(x) + u_i \varphi_i(x) + u_{i+1} \varphi_{i+1}(x).$$

The integral (78) now takes the simplified form

$$\int_{x_{i-1}}^{x_{i+1}} f(u_{i-1}\varphi_{i-1}(x) + u_i\varphi_i(x) + u_{i+1}\varphi_{i+1}(x))\varphi_i(x) dx.$$

plit this integral in two integrals over cell L (left),  $[x_{i-1}, x_i]$ , and cell R (right),  $[x_i, x_{i+1}]$ . Over cell L,  $u$  simplifies to  $u_{i-1}\varphi_{i-1} + u_i\varphi_i$  (since  $\varphi_{i+1} = 0$  on this cell), and over cell R,  $u$  simplifies to  $u_i\varphi_i + u_{i+1}\varphi_{i+1}$ . Make a `sympy` program that can compute the integral and write it out as a difference equation. Give the  $f(u)$  formula on the command line. Try out  $f(u) = u^2, \sin u, \exp u$ .

**hint.** Introduce symbols `u_i`, `u_im1`, and `u_ip1` for  $u_i$ ,  $u_{i-1}$ , and  $u_{i+1}$ , respectively, and similar symbols for  $x_i$ ,  $x_{i-1}$ , and  $x_{i+1}$ . Find formulas for the basis functions on each of the two cells, make expressions for  $u$  on the two cells, integrate over each cell, expand the answer and simplify. You can make L<sup>A</sup>T<sub>E</sub>X code and render it via <http://latex.codecogs.com>. Here are some appropriate python statements for the latter purpose:

```
from sympy import *
...
# expr_i holds the integral as a sympy expression
latex_code = latex(expr_i, mode='plain')
# Replace u_im1 sympy symbol name by latex symbol u_{i-1}
latex_code = latex_code.replace('im1', '{i-1}')
# Replace u_ip1 sympy symbol name by latex symbol u_{i+1}
latex_code = latex_code.replace('ip1', '{i+1}')
# Escape (quote) latex_code so it can be sent as HTML text
import cgi
html_code = cgi.escape(latex_code)
# Make a file with HTML code for displaying the LaTeX formula
f = open('tmp.html', 'w')
# Include an image that can be clicked on to yield a new
# page with an interactive editor and display area where the
# formula can be further edited
text = """
<a href="http://www.codecogs.com/eqnedit.php?latex=%(html_code)s"
target="_blank">

</a>
""" % vars()
f.write(text)
f.close()
```

The formula is displayed by loading `tmp.html` into a web browser.  
Filename: `fu_fem_int.py`.

## Problem 12: Finite elements for the 1D Bratu problem

Address the same 1D Bratu problem as described in Problem 10.

a) Discretize (12) by a finite element method using a uniform mesh with P1 elements. Use a group finite element method for the  $e^u$  term.

b) Set up the nonlinear equations  $F_i(u_0, u_1, \dots, u_{N_x}) = 0$  from a). Compute the associated Jacobian.

Filename: `nonlin_1D_Bratu_fe.pdf`.

## Problem 13: Derive the Newton system from a variational form

We study the multi-dimensional heat conduction PDE

$$\rho c(T) T_t = \nabla \cdot (k(T) \nabla T)$$

in a spatial domain  $\Omega$ , with a nonlinear Robin boundary condition

$$-k(T) \frac{\partial T}{\partial n} = h(T)(T - T_s(t)),$$

at the boundary  $\partial\Omega$ . The primary unknown is the temperature  $T$ ,  $\rho$  is the density of the solid material,  $c(T)$  is the heat capacity,  $k(T)$  is the heat conductivity,  $h$  is a heat transfer coefficient, and  $T_s(T)$  is a possibly time-dependent temperature of the surroundings.

a) Use a Backward Euler or Crank-Nicolson time discretization and derive the variational form for the spatial problem to be solved at each time level.

b) Define a Picard iteration method from the variational form at a time level.

c) Derive expressions for the matrix and the right-hand side of the linear system that arises from applying Newton's method to the variational form at a time level.

d) Apply the Backward Euler or Crank-Nicolson scheme in time first. Then apply a Newton method at the PDE level. Make a variational form of the PDE at a time level.

Filename: `nonlin_heat_Newton.pdf`.

## Problem 14: Derive algebraic equations for nonlinear heat conduction

Consider a 1D heat conduction PDE

$$\rho c(T) T_t = (k(T) T_x)_x,$$

where  $\rho$  is the density of the solid material,  $c(T)$  is the heat capacity,  $T$  is the temperature, and  $k(T)$  is the heat conduction coefficient.

Use a uniform finite element mesh, P1 elements, and the group finite element method to derive the algebraic equations arising from the heat conduction PDE.

a) Discretize the PDE by a finite difference method. Use either a Backward Euler or Crank-Nicolson scheme in time.



) Derive the matrix and right-hand side of a Newton method applied to the discretized PDE.

filename: `nonlin_1D_heat_PDE.pdf`.

### Problem 15: Investigate a 1D problem with a continuation method

Flow of a pseudo-plastic power-law fluid between two flat plates can be modeled by

$$\frac{d}{dx} \left( \mu_0 \left| \frac{du}{dx} \right|^{n-1} \frac{du}{dx} \right) = -\beta, \quad u'(0) = 0, \quad u(H) = 0,$$

where  $\beta > 0$  and  $\mu_0 > 0$  are constants. A target value of  $n$  may be  $n = 0.2$ .

) Formulate a Picard iteration method directly for the differential equation problem.

) Perform a finite difference discretization of the problem in each Picard iteration. Implement a solver that can compute  $u$  on a mesh. Verify that the solver gives an exact solution for  $n = 1$  on a uniform mesh regardless of the cell size.

) Given a sequence of decreasing  $n$  values, solve the problem for each  $n$  using the solution for the previous  $n$  as initial guess for the Picard iteration. This is called a continuation method. Experiment with  $n = (1, 0.6, 0.2)$  and  $n = (1, 0.9, 0.8, \dots, 0.2)$  and make a table of the number of Picard iterations versus  $n$ .

) Derive a Newton method at the differential equation level and discretize the resulting linear equations in each Newton iteration with the finite difference method.

) Investigate if Newton's method has better convergence properties than Picard iteration, both in combination with a continuation method.

## References

- [K95] C. T. Kelley. *Iterative Methods for Linear and Nonlinear Equations*. SIAM, 1995.
- [MLW11] M. Mortensen, H. P. Langtangen, and G. N. Wells. A FEniCS-based programming framework for modeling turbulent flow by the Reynolds-averaged Navier-Stokes equations. *Advances in Water Resources*, 34(9), 2011.

## Index

continuation method, 42, 49

fixed-point iteration, 5

group finite element method, 31

latex.codecogs.com web site, 46

linearization, 5

explicit time integration, 3

fixed-point iteration, 5

Picard iteration, 5

successive substitutions, 5

online rendering of L<sup>A</sup>T<sub>E</sub>X formulas, 46

Picard iteration, 5

product approximation technique, 31

relaxation (nonlinear equations), 10

single Picard iteration technique, 6

stopping criteria (nonlinear problems),  
6, 18

successive substitutions, 5